

(19)



JAPANESE PATENT OFFICE

PATENT ABSTRACTS OF JAPAN

(11) Publication number: **10207891 A**

(43) Date of publication of application: **07 . 08 . 98**

(51) Int. Cl. **G06F 17/27**
G06F 17/30

(21) Application number: **09006777**

(71) Applicant: **FUJITSU LTD**

(22) Date of filing: **17 . 01 . 97**

(72) Inventor: **NAKAO YOSHIO**

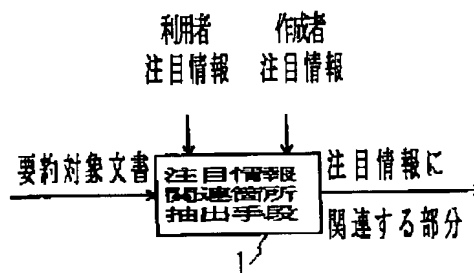
(54) **DOCUMENT SUMMARIZING DEVICE AND ITS METHOD**

COPYRIGHT: (C)1998,JPO

(57) Abstract:

PROBLEM TO BE SOLVED: To easily grasp how to information noticed by a user is treated in a document only by reading out its summary by extracting parts related to the user's noticed information and writer's noticed information from a document to be summarized based on both the information.

SOLUTION: An extraction means 1 extracts a part related to noticed information as the center part of a summary based on the user's noticed information and writer's noticed information. In this case, an information notice reference for judging that a part including more noticed information is more important is used and the center part of the summary is determined. When two kinds of the noticed information are considered, a summary adopting both of the information requested by the user and the important information in the document, i.e., contents to be written by the writer, is prepared. When either one of the two sorts of the noticed information is weighted, a summary attaching greater importance to the user or the writer is prepared.



(19)日本国特許庁 (J P)

(12) 公 開 特 許 公 報 (A)

(11)特許出願公開番号

特開平10-207891

(43)公開日 平成10年(1998) 8月7日

(51)Int.Cl.⁸

識別記号

F I

G 0 6 F 17/27
17/30

G 0 6 F 15/20 5 5 0 A
15/401 3 2 0 A
15/403 3 3 0 C

審査請求 未請求 請求項の数23 O L (全 28 頁)

(21)出願番号

特願平9-6777

(22)出願日

平成9年(1997) 1月17日

(71)出願人 000005223

富士通株式会社

神奈川県川崎市中原区上小田中4丁目1番
1号

(72)発明者 仲尾 由雄

神奈川県川崎市中原区上小田中4丁目1番
1号 富士通株式会社内

(74)代理人 弁理士 大菅 義之 (外1名)

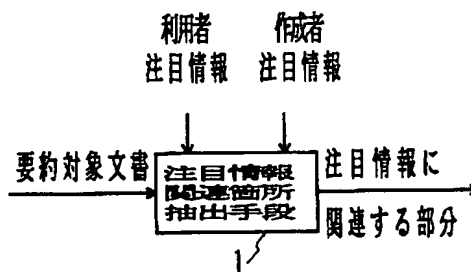
(54)【発明の名称】 文書要約装置およびその方法

(57)【要約】

【課題】 自然言語による、電子化された文書の要約を行う装置に関し、利用者の注目していることと、利用者が既に持っている知識に応じた要約作成を可能とする。

【解決手段】 要約対象文書の中で、要約の利用者が注目している情報としての利用者注目情報と、文書の作成者が注目を促している情報としての作成者注目情報とに基づいて、これら2つの注目情報に関連する部分を抽出する手段1を備え、その抽出結果に応じて要約を作成する。

本発明の第1の原理の説明図



【特許請求の範囲】

【請求項1】 計算機可読の文書の選択・閲覧・編集および管理の支援のために文書要約を行う装置において、要約対象文書の内容の中で、要約を利用する利用者が注目している情報としての利用者注目情報、および該要約対象文書の作成者が注目を促している情報としての作成者注目情報に基づいて、該要約対象文書中の該2種類の注目情報に関連する部分を抽出する注目情報関連箇所抽出手段を備えることを特徴とする文書要約装置。

【請求項2】 前記利用者注目情報が、前記要約対象文書の検索のために利用者から入力される質問文の内容であることを特徴とする請求項1記載の文書要約装置。

【請求項3】 前記利用者注目情報および／あるいは作成者注目情報が単語列あるいは重みづけられた単語列の形式であることと、

前記注目情報関連箇所抽出手段が、前記要約対象文書における該単語列内の単語の出現の程度に応じて、前記2種類の注目情報に関連する部分を抽出することを特徴とする請求項1記載の文書要約装置。

【請求項4】 前記文書要約装置において、前記利用者が興味を有する事柄を利用者嗜好特性としてあらかじめ蓄積する利用者嗜好特性蓄積手段を更に備え、

前記注目情報関連箇所抽出手段が、該利用者嗜好特性蓄積手段の蓄積内容を前記利用者注目情報として利用することを特徴とする請求項1記載の文書要約装置。

【請求項5】 前記文書要約装置において、前記利用者嗜好特性蓄積手段が複数の利用者のそれぞれに対して利用者嗜好特性を蓄積すると共に、あらかじめ定められたアクセス制御方式のもとで、前記要約を利用する利用者の利用者注目情報として、異なる利用者の利用者嗜好特性を含む情報を前記注目情報関連箇所抽出手段に与えて前記2種類の注目情報に関連する部分を抽出させる他利用者嗜好特性活用手段を更に備えることを特徴とする請求項4記載の文書要約装置。

【請求項6】 前記作成者注目情報が、通常の流通文書に含まれ、かつ作成者が文書の要点をまとめて提示している情報であって、文書の表題、文書中の章・節および図表の見出し、目次、用語・事項の索引の情報であることを特徴とする請求項1記載の文書要約装置。

【請求項7】 前記文書要約装置において、複数の要約対象文書に対するそれぞれの作成者注目情報をマージする作成者注目情報マージ手段を更に備え、マージした作成者注目情報に基づいて、前記注目情報関連箇所抽出手段が該複数の要約対象文書中の前記2種類の注目情報に関連する部分を抽出して、該複数の要約対象文書間の比較情報とすることを特徴とする請求項1記載の文書要約装置。

【請求項8】 前記文書要約装置において、文書作成時点以後に該文書の作成者あるいは文書管理者

によって指定される作成者注目情報を、該作成者注目情報に対応する文書と共に格納する文書格納手段を更に備え、

前記注目情報関連箇所抽出手段が、該文書格納手段に格納されている作成者注目情報を利用することを特徴とする請求項1記載の文書要約装置。

【請求項9】 計算機可読の文書の選別・閲覧・編集および管理の支援のために文書要約を行う装置において、利用者がすでに知っている利用者既知情報、および／あるいは該要約作成時点において、過去に利用者に提示された文書に基づいて利用者がすでに知っているものとみなせる履歴的既知情報と、該2種類の既知情報以外の情報を区別して使用して要約を作成し、要約の可読性を向上させる要約可読性向上手段を備えることを特徴とする文書要約装置。

【請求項10】 前記利用者既知情報、および／あるいは履歴的既知情報が既知概念と既知の事柄とによって構成されることと、前記要約可読性向上手段が、要約内の既知でない概念を減少させ、かつ既知でない事柄については既知性の低いものを優先して要約に取り入れるることによって要約の可読性を向上させることを特徴とする請求項9記載の文書要約装置。

【請求項11】 前記文書要約装置において、前記概念の既知性が文書内に出現する用語の既知性であることと、文書内に出てくる用語を認定する用語認定手段と、該用語認定手段によって認定された用語の既知性を判定する用語既知性判定手段とを更に備えることを特徴とする請求項10記載の文書要約装置。

【請求項12】 前記文書要約装置において、前記事柄の既知性が文書内に出現する用語の組み合わせの既知性であることと、文書内に出てくる用語の組み合わせを認定する用語組み合わせ認定手段と、該用語組み合わせ認定手段によって認定された用語の組み合わせの既知性を判定する用語組み合わせ既知性判定手段とを更に備えることを特徴とする請求項10記載の文書要約装置。

【請求項13】 前記文書要約装置において、前記事柄の既知性が文書内に出てくる用語と述語との組み合わせの既知性であることと、文書内に出てくる用語と述語との組み合わせを認定する用語と述語の組み合わせ認定手段と、該用語と述語の組み合わせ認定手段によって認定された用語と述語の組み合わせの既知性を判定する用語と述語の組み合わせ既知性判定手段とを更に備えることを特徴とする請求項10記載の文書要約装置。

【請求項14】 前記文書要約装置において、前記利用者が熟知している事柄を利用者知識としてあら

はじめ蓄積する利用者知識蓄積手段を更に備え、前記要約可読性向上手段が、該利用者知識蓄積手段に蓄積されている利用者知識を前記利用者既知情報として利用することを特徴とする請求項9記載の文書要約装置。

【請求項15】 前記文書要約装置において、前記利用者知識蓄積手段が複数の利用者のそれぞれに対する利用者知識を蓄積すると共に、

あらかじめ定められたアクセス制御方式のもとで、前記要約を利用する利用者の利用者既知情報として、異なる利用者の利用者知識を含む情報を前記要約可読性向上手段に用いさせる他利用者知識活用手段を更に備えることを特徴とする請求項14記載の文書要約装置。

【請求項16】 前記文書要約装置において、該文書要約装置あるいは該文書要約装置を含むシステムの稼働期間において利用者に提示された文書や要約を利用者の閲覧履歴として保持し、該閲覧履歴を前記履歴的既知情報の基として前記要約可読性向上手段に与える閲覧履歴保持手段と、

該閲覧履歴保持手段に保持されている文書や要約と要約対象文書とを相互参照する文書相互参照手段とを更に備えることを特徴とする請求項9記載の文書要約装置。

【請求項17】 前記閲覧履歴保持手段が、前記稼働期間を含む長期間に渡る複数の利用者の閲覧履歴をそれぞれの利用者毎に保持することを特徴とする請求16記載の文書要約装置。

【請求項18】 前記文書要約装置において、あらかじめ定められたアクセス制御方式のもとで前記要約を利用する利用者の履歴的既知情報として、異なる利用者の閲覧履歴に基づく履歴的既知情報を含む情報を前記要約可読性向上手段に与える他利用者閲覧履歴活用手段を更に備えることを特徴とする請求項17記載の文書要約装置。

【請求項19】 前記文書要約装置において、要約対象文書の中の各文を、文の述語と該述語に支配される名詞を基本として構成される述語句に分割し、該述語句のうちで他の述語句に依存していない述語句を主述語句とし、該述語句が主題句を含むときは該主題句を分離し、1つの文内または他の文の間での構文的依存構造に従って主題句と主述語句、主述語句と他の述語句との間に依存関係の設定を行い、該設定結果をともなう文書内容を前記要約可読性向上手段に与える文分割・依存関係設定手段を更に備えることを特徴とする請求項9記載の文書要約装置。

【請求項20】 計算機可読の文書の選別・閲覧・編集および管理の支援のために文書要約を行う装置において、

要約対象文書の内容の中で、要約を利用する利用者が注目している情報としての利用者注目情報、および該要約対象文書の作成者が注目を促している情報としての作成者注目情報に基づいて、該要約対象文書中の該2種類の

注目情報に関連する部分を抽出する注目情報関連箇所抽出手段と、

該抽出結果に対して、利用者がすでに知っている利用者既知情報、および／あるいは該要約作成時点において、過去に利用者に提示された文書に基づいて利用者がすでに知っているときとみなせる履歴的既知情報と、該2種類の既知情報以外の情報を区別して使用して要約を作成し、要約の可読性を向上させる要約可読性向上手段とを備えることを特徴とする文書要約装置。

10 【請求項21】 計算機可読の文書の選別・閲覧・編集および管理の支援のために文書要約を行う方法において、

要約対象文書の内容の中で、要約を利用する利用者が注目している情報としての利用者注目情報、および該要約対象文書の作成者が注目を促している情報としての作成者注目情報に基づいて、該要約対象文書中の該2種類の注目情報に関連する部分を抽出することを特徴とする文書要約方法。

20 【請求項22】 計算機可読の文書の選別・閲覧・編集および管理の支援のために文書要約を行う方法において、

利用者がすでに知っている利用者既知情報、および／あるいは該要約作成時点において、過去に利用者に提示された文書に基づいて利用者がすでに知っているときとみなせる履歴的既知情報と、該2種類の既知情報以外の情報を区別して使用して要約を作成し、要約の可読性を向上させることを特徴とする文書要約方法。

30 【請求項23】 計算機可読の文書の選択・閲覧・編集および管理の支援のために文書要約を行う方法において、

要約対象文書の内容の中で、要約を利用する利用者が注目している情報としての利用者注目情報、および該要約対象文書の作成者が注目を促している情報としての作成者注目情報に基づいて、該要約対象文書中の該2種類の注目情報に関連する部分を抽出し、該抽出結果に対して、利用者がすでに知っている利用者既知情報、および／あるいは該要約作成時点において、過去に利用者に提示された文書に基づいて利用者がすでに知っているときとみなせる履歴的既知情報と、該2種類の既知情報以外の情報を区別して使用して要約を作成し、要約の可読性を向上させることを特徴とする文書要約方法。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、電子化された自然言語による文書の要約を行う装置に関するものであり、主として、検索された大量の文書の選別や閲覧、あるいは、蓄積された大量の文書の閲覧、再構成（再利用）や管理のプロセスを支援するために、使われることを意図したものである。

【0002】近年、文書の電子化が進み、大量の文書を

蓄積して再利用することで文書作成のコストを大幅に引き下げることが可能になってきた。また、一方で、技術の高度化に伴うマニュアル類のボリュームの増加と多様化、インターネットに代表される新たな文書流通メディアの出現もあいまって、計算機上で閲覧や再利用のための再構成操作を行える文書の量が爆発的に増加している。

【0003】このような大量の文書を利用するためには、まず、個々の文書の有用性を素早く判定し、利用目的にあった文書を選別することが重要である。そのためには、文書一覧に文書内容をイメージできるような情報を合わせて提示することが求められる。このような目的にあった情報としては、文書の見出しや抄録などがあるが、見出しが十分に文書内容を表現していない場合や、抄録がない場合も多い。また、特にオンラインで文書を閲覧する場合には、表示可能な文字数が限られるため、例えば抄録が作成されていても長過ぎて一覧表示に適さないこともある。そこで、適切な要約を自動的に生成する技術が強く求められることになる。

【0004】また、文書の再利用を効率的に行うためには、大量の文書を分類・整理して蓄積することが求められる。この場合にも、新たな文書を分類するために当該文書の内容を素早く把握したり、蓄積文書の管理者が分類体系を改良するために分類状況を概観したり、分類体系に通じていない利用者がどのような分類がなされているのかを把握したりすることなどを支援する意味で、やはり適切な要約が望まれる。

【0005】本発明は、このような用途をもつ文書要約装置において、利用者が何に注目しているのか、また、利用者がどのような知識を既にもっているのかに応じて、要約結果を調整する点に特徴をもつものである。

【0006】

【従来の技術】従来の文書の要約作成の技術には、大きく分けて2つの方法がある。第1の方法は、文書において重要な部分（通常は、文、段落、節などの文書の論理要素。以下「文」で代表させる。）を認定し、それを抽出することで要約を作成するものである。第2の方法は、要約として抽出すべき情報の型紙を用意して、その型紙の条件にあった文書中の語句を抽出して要約としたり、その型紙によくあてはまる文を抽出して要約とする方法である。第2の方法は本発明との関連性は低いので、ここでは第1の方法を説明する。

【0007】第1の方法は、さらに、何を手がかりに文の重要性を評価するかによっていくつかの方法に分類される。代表的な方法としては、①文書中に出現する単語の頻度と分布を手がかりとするもの、②文と文とのつながり方や文の出現位置を手がかりとするものの2つがある（その他、文の構文的パターンによって重要性を評価するものなどもあるが、本発明との関連性が低いので省略する）。

【0008】①の方法、すなわち、文書中に出現する単語の頻度と分布を手がかりとするものは、通常、まず文書中に含まれる単語（語句）の重要度を決定し、次に重要な単語をどれ位含んでいるかで文の重要度を評価し、重要な文を選択するという手順で要約を作成する。単語の重要度を決定する手法としては、ある文書内の出現度数そのままを用いたり、あるいは、一般的な文書集合における単語の出現度数とのずれなどを加味して重みをつけたり、あるいは単語の出現位置に応じて重みをつけたり（見出しに出現する語を重要とみなすなど）する方法が知られている。なお、対象とする単語は、日本語であれば自立語（特に名詞）、英語であれば内容語のみに限るのが通例である。ここで自立語・内容語とは実質的な意味を持つ名詞、形容詞、動詞などの語であり、助詞や前置詞、形式名詞など、専ら構文的役割を示すために使われる語と区別されるものである。なお、日本語の自立語の形式的定義は、独立した文節を構成できる語というもので、前記の説明とは若干のずれがあるが、対象とする単語を自立語に限ることの趣旨は前記の説明中の区別による。

【0009】このようなものには、例えば以下のようなものがある。特開平6-259424「文書表示装置及び文書要約装置並びにディジタル複写装置」およびその発明者による次の文献1では、見出しに含まれる単語を多く含む部分を（見出しに関連の深い）重要な部分として抽出することで要約を作成している。

【0010】文献1：亀田雅之、擬似キーワード相関法による重要キーワードと重要文の抽出、第2回年次大会、pp.97～100、言語処理学会1996年3月、

特開平7-36896「文書を要約する方法および装置」では、文書中に現れる表現（単語など）の複雑さ（語の長さなど）から重要な表現の候補（シード）を選び、重要性の高いシードをより多く含む文を抽出することで要約を作成している。

【0011】特開平成8-297677「主題の要約を生成する自動的な方法」では、文書内の単語の出現頻度が大きい順に「主題の用語」を認定し、重要な「主題の用語」をより多く含む文を抽出することで要約を作成している。

【0012】特開平成6-215049「文書要約装置」では、検索結果と質問文との関連性判定などによく用いられるベクトル空間モデルを適用して、文書全体の特徴ベクトルと最も類似した特徴ベクトルを持つ文や段落から文を選択していくことで要約を作成している。ここで、ベクトル空間モデルとは、キーとなる単語ごとに、あるいは単語の意味成分ごとに、次元（軸）をわりあて、文書や質問文におけるそれらの単語の出現の有無あるいは出現頻度の値の列（特徴ベクトル）で、文書や質問文の特徴を表現するものである。

【0013】②の方法、すなわち、文と文とのつながり

方や文の出現位置を手がかりとするものとは、順接・逆接、展開などの文の接続関係（結束関係と呼ぶ場合もある）や、文が出現している文書中の位置などをもとに、文の（相対的）重要性を判定し、重要な文を選択するものである。例えば、特開平7-182373「文書情報検索装置及び文書検索結果表示方法」およびその発明者らによる次の文献2、別の著者による文献3がある。

【0014】文献2：住田一男、知野哲朗、小野頭司、三池誠司、文書構造解析に基づく自動抄録生成と検索提示機能としての評価・電子情報通信学会論文誌、Vol. J78-D-II, No. 3, pp. 511~519, 1995年3月。

文献3：山本和英、増山繁、内藤昭三、文書内構造を複合的に利用した論説文要約システムGREEN。情処研報NL-99-3, 情報処理学会1994年1月。

以上のように文書全体の要約を作成する技術の他に、個々の文書の有用性の判定を支援するための技術として、利用者の注目している部分を提示する技術もある。周知の技術として、KWIC (Keyword In Context) と呼ばれる検索語の周囲を表示する方法や、それに類似した検索語の近傍表示の方法が広く使われている。

【0015】また、論文において研究の背景となる事情を述べた部分だけを提示したり、新聞の第一パラグラフだけを提示するなど、利用者の目的に応じて特定の部分だけを提示する方法もある。例えば、前掲の特開平成7-182373や文献3、別の著者による次の文献4、5がある。ただし、これらは、分野に特有な文書構成の類型や言い回しなどを手がかりとして、文書の論理構造上で特別な機能を持つ部分を選択するものであり、利用者が注目している内容に着目して、その内容と関連性の高い部分を提示しているわけではない。

【0016】文献4：神門典子。複数領域における日本語原著論文の機能構造分析：構成要素カテゴリの自動付与。Library and Information Science, No. 31, pp. 25~38, 1993年。

文献5：神門典子、原著論文の機能構造の分析とその応用。図書館学会年報、Vol. 40, No. 2, pp. 49~61, 1994年6月。

要約の可読性を低下させる要因としては、冗長な表現、利用者の知らない用語（未知の用語）の出現、解消されていない照応表現（anaphoric expression：日本語ならいわゆるコソアド語など）の出現などがある。

【0017】このうち、冗長な表現については、修飾要素と被修飾要素の語彙的な性質や関係、修飾要素と被修飾要素の距離などによるヒューリスティックにより、余分な修飾要素を削除する手法が知られている。例えば、前掲の文献3では、日本語の新聞記事中の文において同じ名詞に対して2つ以上の修飾要素があった場合、前の修飾要素を削除するというヒューリスティックが示され、また、同じ著者らによる次の文献6では、同一事件に関する一連の記事で、特有の言い回しから導入部と認

定された部分で出現する名詞の7割以上が以前の記事で既出の場合にその導入部を削除するというヒューリスティックが示されている。

【0018】文献6：船坂貴浩、山本和英、増山繁。冗長度削減による関連記事の要約。情処研報NL-114-7, 情報処理学会, 1996年7月。

未知の用語の出現については、用語の定義や説明をしている部分が文書中にあれば、それを要約に含めればよいことは自明である。このような部分を探すには、初出の部分あるいは用語の定義を示すマーク（日本語なら「とは」など）の付いた当該用語が出現する箇所を探せばよいことも、学校教育レベルの知識である。

【0019】照応表現についても、その先行詞（antecedent）を探し、照応表現を先行詞で置き換えたり、先行詞を含む部分を要約に含めれば、要約が理解しやすいものとなることは自明である。照応表現の先行詞の同定に関しては、センタリング（centering）と呼ばれる手法などが知られている。これは、後続の文で照応表現になりやすい要素（センタ）を構文的役割などに応じて優先度つきで認定しておき、後続の照応表現の現れ方による制約を加味して、先行詞をセンタの中から選択するというものである。なお、類似の手法で、センタと同様な概念を焦点（focus）と呼ぶものもある。ただし、いずれにしても完全な結果を得ることができる技術ではない。センタリングの手法については次の文献7、8がある。

【0020】文献7：Megumi Kameyama. A property-sharing constraint in centering. In Proceedings of the 24th Annual Meeting of Association for Computational Linguistics, pp. 200~206, 1986年。

文献8：Susan E. Brennan, Marilyn W. Friedman, and Carl J. Pollard. A centering approach to pronouns. In Proceedings of the 25th Annual Meeting of Association for Computational Linguistics, pp. 155~162, 1987年。

前掲の特開平7-182373や、同一発明者らによる特開平7-44566「抄録作成装置」では、このような手法を実装して用語の定義箇所や照応表現の先行詞を推定し、元の用語や照応表現からハイパーテキストチュアル（hyper-textual）なリンクを設定するなどして利用者の便を図っている。

【0021】

【発明が解決しようとする課題】大量の文書から有用な文書を選別するためには、利用者にとっての文書の有用性を素早く判定できるよう、利用者の求める情報を文書の作成者がどのように取り扱っているのかを示すことが重要である。検索システムでは、利用者が得たいと思っている情報は、質問文やキーワードによる検索式として表現されることが多い。しかしながら、質問文や検索式中の単語が検索対象の文書の中に見つかったからといって、そこに利用者の求める情報が書かれているとは限ら

ない。たとえば、特許公報を「翻訳」というキーワードで検索すると、利用者の得たい情報は「自然言語の文の翻訳」に関する特許情報なのに、検索結果には「機械語の翻訳」に関する情報が大量の特許が含まれることがある。この場合、「翻訳」という語がどのような文脈で使われているのかを提示すれば、ある程度、文書の選別を支援することが可能である。このような目的で、前節に挙げたKWICを用いることもできるが、物理的な近傍を表示するだけなので論旨の流れなどを把握することが難しく、簡潔で要を得た要約を提示することは必ずしもできない。

【0022】この観点からみると、前述のように、従来の要約作成の技術では専ら文書中の文の重要度のみを要約中に文を含めるかどうかの基準としており、利用者の注意がいずれに向けられているのかが考慮されていない、という問題がある。そのため、例えば検索システムに入力したキーワードが、言語学の文献の用例の部分で照合されて取り出されてしまった場合などでは、自動作成された要約中にはキーワードが含まれず、利用者の入力とどのような関連がある文献なのかを判定しづらくなる場合などが考えられる。

【0023】言語学の文献に関連する問題点についてさらに説明する。言語学の文献は、言語の形式的な性質を議論するもので、そこで取り上げられている用例の内容については言語学と関係している必要はない。例えば、

「象は鼻が長い」や「僕はうなぎだ（料理の注文の意）」というのは国語学では有名な例文である。動物のことを調べようとしている時に、このような例文で出ている言語学の文献が検索されてしまう可能性がある。言語学の文献のために、文献中に含まれる語彙の頻度分布などをとれば、動物に関する語彙は少なく、例えば

「象」はあまり重要でない語彙であると計算されてしまう。そうすると、検索結果の表示に語彙の頻度分布などに基づいて自動生成した要約を用いる場合、このような例文は要約中に含まれにくくなり、不都合が生じる可能性がある。つまり、「象」というキーワードを入力した時に、このような言語学の文献が検索され、それに関わらず、検索結果の表示（自動生成した要約）の中には「象」が含まれず、何でこのようなものが検索されたのかが理解できなくなる可能性があるということになる。逆に、キーワードの近傍表示だけしかない場合には、用例の部分だけが表示され、どういう趣旨の文献なのかを理解できなくなることもある。

【0024】もう一つの問題点として、従来の要約作成の技術には、利用者の知識レベルに合わせて要約を作成する手段を備えていない、という問題点もある。利用者の知識レベルは、利用者ごとに大きく異なる可能性があり、特に専門的な用語に関する知識がどの程度あるかに合わせて要約に用語に関する定義や説明の記述部分を含めるかどうかを切り換えないと、知識レベルの高い利用

者にとっては冗長な要約になったり、知識レベルの低い利用者にとっては理解し難い要約になってしまったりすることになる。

【0025】以上のことから、本発明は要約作成において、利用者の注意の方向が考慮されていない、利用者の知識レベルが考慮されていない、という2つの問題点を解決し、それらを統一的に扱う手段を提供することを目的とする。

【0026】

10 【課題を解決するための手段】図1、および図2は本発明の原理説明図である。これらの図は、自然言語による電子化された文書の選別・閲覧・編集、および管理の支援のために、文書要約を行う文書要約装置の原理を説明するものである。

【0027】図1は本発明の第1の原理の説明図である。同図において注目情報関連箇所抽出手段1は、利用者の注目情報と作成者注目情報とに基づいて、要約対象文書中でこれら2つの情報に関連する部分を抽出するものである。ここで利用者注目情報とは、要約対象文書の内容の中でその要約を利用する利用者が注目している情報であり、また作成者注目情報とは要約対象文書の作成者が利用者に対して注目を促している情報である。

20 【0028】図2は本発明の第2の原理の説明図である。同図において要約可読性向上手段2は、利用者既知情報と履歴的既知情報とに基づいて、これら二種類の既知情報以外の情報と二種の既知情報とを区別して用いて要約を作成し、要約の可読性を向上させるものである。ここで利用者既知情報とは、要約対象文書の内容の中で、その要約を利用する利用者がすでに知っている情報であり、また履歴的既知情報とは、要約作成時点において、それ以前に利用者に提示された文書に基づいて利用者がすでに知っていると思わせる情報である。

30 【0029】後述する本発明の実施例においては、図1と図2とによって説明される2つの原理を同時に用いて、要約の作成が行われる。まず図1においては、利用者注目情報および作成者注目情報という2種類の注目情報に基づいて、注目情報に関連する部分が要約の中心部分として抽出されることになる。ここではこれらの注目情報を多く含む部分ほど重要であると判定する情報の注目性基準が用いられて、要約の中心部分が決定される。

40 【0030】2種類の注目情報を考慮することによって、利用者が求める情報と、文書において重要な情報、すなわち作成者が書こうとしていた内容の双方を取り込んだ要約を作成することができる。これら2種類の注目情報のいずれかに重みをつけることによって、利用者注目情報だけを重視した要約から、文書における重要性だけを重視した要約まで、目的に応じた要約を作成することができる。あるいはこれら2種類の注目情報を同等に扱うことによって、利用者が求める情報と、作成者が書こうとしていた内容との双方をバランスよく抽出した要

約を作成することもできる。

【0031】図2においては、利用者既知情報と履歴的既知情報との2つに基づいて要約が作成される。これは利用者の知識レベルに合わせて要約を作成することを意味し、これによって要約の可読性を向上させることができる。この可読性の向上のために情報の既知性基準が用いられる。情報の既知性基準は、例えば概念の既知性基準と事柄の既知性基準との2つの基準を意味する。

【0032】概念の既知性基準とは、要約を構成する要素概念、特に主題に関する要素概念が原則として既知でなくてはならないという基準である。ここで要素概念とは、要約に含まれる個々の語句が表す概念のことである。言い換えれば、要約に出力する用語（主として名詞）は、原則として利用者にとって既知でなくてはならないという基準である。本発明においては、この基準に基づいて、利用者に理解できない用語に関しては、例えば必要な説明を追加して要約が作成される。

【0033】事柄の既知性基準とは、例えば文書の中に出現する用語の組み合わせに関するものであり、その組み合わせ全体で述べられている事柄（事実、あるいは命題）については、既知性の低いものほど優先されて要約に取り入れられる。

【0034】この事柄の既知性基準によって通常の場合、すなわち独立した文書1つを要約する場合などは、要約の中で同じ事柄が何度も出力されるのが抑制されることになる。また関連する複数の文書を一括して要約する場合、特に同一事件に関する一連の記事や、記載コラムなどをまとめて、要約対象の文書の間の関係を明らかにして提示できる場合には、同一の事柄に関する記述を削減するための基準として用いられる。

【0035】この事柄の既知性基準は、前述の情報の注目性基準によって注目情報の含み方が同程度と判定された文が複数あった場合には、事柄の既知性の低い方を選んでその低い方の文の内容が要約に含められるという意味で、弱めの制約となるものである。すなわち、概念の既知性基準が「（原則として）既知でなければならない」のに対して、事柄の既知性基準は「既知性の低いものが優先される」だけであるために弱めの制約となる。

【0036】このように本発明においては、情報の注目性基準と情報の既知性基準という2つの基準を用いて、要約の作成が行われる。

【0037】

【発明の実施の形態】図3は本発明の文書要約装置の構成を示すブロック図である。同図において、文書要約装置は要約プロセス制御部10、文書構造解析部11、文解析部12、文分割・依存関係設定部13、文選択部14、要約整形部15を基本構成要素として備えている。このうち本発明にとって特徴的な構成要素は、文分割・依存関係設定部13と、文選択部14である。

【0038】要約プロセス制御部10は、利用者と装置

とのインタフェースとなると共に、文書要約装置の動作全体を制御するものである。利用者との間のインタフェースとしては、利用者が注目している情報、すなわち利用者注目情報や、要約作成に関する要望などの入力を受け取り、文書要約のプロセスを適切に起動し、要約結果を利用者に出力することになる。利用者注目情報の代表的な形式は利用者から入力される質問文であるが、求める情報に関するキーワードや、読書案内に載っている紹介文の形式とすることも可能である。

10 【0039】要約作成に関する要望、すなわち要約作成に関する制約情報としては、利用者から必須出力要素が指定されると共に、その他の制御命令が与えられる。ここで必須出力要素とは、例えば見出しのように、要約の中に必ず含めるべき要素である。その他の制御命令としては、注目情報や既知情報としてどのような情報を利用するか、またそれらの情報をどのように利用するか、望ましい要約の長さ、要約処理でどのような単位を基本として要約を作成するかなどがある。この基本単位としては、通常は文、あるいは述語句が用いられる。

20 【0040】また、要約プロセス制御部10では、複数文書の比較を支援するために本装置を用いる場合に、まず個々の比較対象文書に関する文書構造認識解析および文解析までの処理を行い、その出力を集計して作成者注目情報のマージし、マージした注目情報にもとづいてそれ以降の要約処理（文分解・依存関係設定処理、文選択処理、要約成形処理）を行うよう、各処理部の動作の制御および処理経過の記憶も行う。

30 【0041】図3において、文書要約装置のメモリには、利用者の嗜好特性16、利用者の知識17、および閲覧履歴18が蓄積されると共に、例えば他のメモリに入力文書（群）19が格納される。

【0042】利用者の嗜好特性16は、利用者が興味を持っている事柄を蓄積するものである。ここには利用者が、例えば自己紹介の際に使うような趣味の説明文や、利用者が興味を持った文書そのままを蓄積したり、そのような文書の中から出現頻度の大きいキーワードを抽出して保存したり、利用者が検索に際してよく使うキーワードや質問文を保存しておいてもよい。

40 【0043】利用者の知識17は、利用者がよく知っている情報を、利用者既知情報として蓄積するものである。ここには例えば利用者が知っている専門用語のリストなどが蓄積される。

【0044】閲覧履歴18は、利用者がどのような文書や要約をいつ頃閲覧したかという履歴を蓄積するものである。入力文書（群）19、基本的には要約対象文書を格納するものであり、通常はどのような形式の電子化文書でも用いることができる。具体例としては、電子出版等で用いられている文書構造記述言語のSGML（スタンダードジェネライズドマークアップランゲージ。ISO8879）を用いればよい。要約対象文書に対し

て、例えば作成者あるいは文書の管理者によって、文書の作成時点以後に指定された作成者注目情報を、文書と対応させて蓄積することもできる。

【0045】メモリの内容としての利用者の嗜好特性16、すなわち利用者注目情報、利用者の知識17、すなわち利用者既知情報、閲覧履歴18、すなわち履歴的既知情報、入力文書(群)19すなわち要約対象文書は要約プロセス制御部10によって管理され、要約の作成に使用される。

【0046】文書構造解析部11は、要約プロセス制御部10から要約対象文書や必須出力要素の指定内容などを受け取り、文書の構造を解析し、文書内容を文解析部12に出力すると共に、依存関係付文書構造情報を文分割・依存関係設定部13に与えるものである。

【0047】文書構造解析部11は、まず文書の書式やマークアップ情報などから、見出しや本文というような文書の論理的構成要素を認定し、例えば見出しと本文、あるいはそれと同様の関係を持つ要素として列挙構造の項目名とその内容などを対応づけて、構成要素間の依存関係を抽出する。この依存関係では、例えば本文内の要素が従属ブロック、見出し内の要素が依存先とする。

【0048】文書構造解析部11は、作成者注目情報として文書の論理的構成要素、すなわち章、節・図表などの見出しや、目次、用語や事項の索引などを用いる場合に、作成者注目情報にあたる部分を認定し、作成者注目の印をつけて、文解析部12および文分割・依存関係設定部13を介して文選択部14に与える。要約プロセス制御部10によって必須出力要素が指定された場合には、該当する部分に必須出力の印をつけて、同様に文選択部14に与える。文解析部12が1つの文を単位としてしか処理できない場合には、文の認定も行うことになる。

【0049】文書構造解析部11による具体的な処理は、文書の種類、例えば単なる自然言語の文書か、構造化された文書(例えばSGML文書)かなどによって異なり、本発明にとって本質的なものではないので、その詳細の説明は省略する。

【0050】文解析部12は、文書の内容を文書構造解析部11から受け取り、それに含まれる単語を認定し、単語の出現位置や品詞情報をつけた単語列の形で文書内容を文分割・依存構造解析部13に出力するものである。また、利用者注目情報が質問文などの自然言語の形で与えられた場合には、自然言語の利用者情報からも同様に単語列を作成し、文選択部14に出力する。具体的な処理の方法としては、形態素解析法として各種のものが知られており、それを用いればよいので説明は省略する。なお、単語列に付与される出現位置とは、文構造解析部11から出力される依存関係付文書構造と単語列とを対応付けるものであり、文分割・依存構造解析部13では文書構造解析部11で設定された文書の構成要素間

の依存関係を述語句間の関係に変換するために使われ、文選択部14では文構造解析部11で設定された必須出力の印および作成者注目の印に従い、必須に出力する述語句や作成者注目情報に対応する注目概念を認識するために使用される。

【0051】文分割・依存関係設定部13は、文書構造解析部11から出力される依存関係付文書構造情報と、文解析部12から出力される出現位置や品詞情報がつけられた単語列、および要約プロセス制御部10から出力される既知概念を用いて、後述する文分割処理と依存関係設定処理を行い、文選択部14に対して依存関係付述語句列(述語句リスト)を出力するものである。

【0052】文選択部14は、要約プロセス制御部10や文解析部12から出力される注目情報と、要約プロセス制御部10から出力される既知の事柄を示す情報に従って、文分割・依存関係設定部13から出力される依存関係付述語句列に対して後述する文選択処理を実行し、要約に含めるべき重要な述語句(文)を選択して、後述する選択結果リストを作成するものである。

【0053】ここで、質問文などの自然言語の注目情報が文解析部12に与えられ、その他の注目情報が文解析部12を経由することなく、直接文選択部14に与えられる理由を説明する。本実施例においては、後述するように文の重要度としての注目情報量は、単語、例えば名詞を単位に計算される。そこで自然言語文として注目情報が与えられた場合には、その自然言語文を単語に分割する必要がある。質問文や文書から取り出された見出しなどが、文解析部12を経由して文選択部14に渡されるのはそのためである。

【0054】一方、例えば利用者の嗜好特性16として蓄積されている利用者注目情報などは、あらかじめ文解析を行った後に適切な形式でメモリに格納することが可能であり、この場合に文解析部12を経由することなく、文選択部14に直接に与えることができる。なお本実施例では自然言語と無関係な情報は利用しないが、メモリへの格納形式は後述する意味ネットワーク表現であったり、フレーム表現であったりしてもよく、蓄積された情報は単なる自然言語に限られない。

【0055】文選択部14で用いられる注目情報としては、代表的には名詞のリストが与えられる。注目度の高い名詞には注目度に対応する重みを与えることもできる。また名詞以外の自立語(動詞や形容詞など)を注目情報として与えたり、名詞と用言の組として与えることも可能である。以下の説明では、重みなしの名詞リストが注目情報として与えられる場合を中心に実施例を説明する。なお述語句の選択においては、文分割・依存関係設定部13の処理によって設定された概念の既知性基準に適合した述語句(文)間の依存関係に違反しないように選択処理が行われるため、概念の既知性基準にも適合した選択結果リストが作成されることになる。

【0056】要約整形部15は、文選択部14によって選択された文を元の文書における出現順に並べ、必要に応じて抽出されなかった文の存在を表す印や、段落の境界などを挿入し、要約を読みやすい形式に整形する。履歴的既知情報への依存関係が設定されている場合には、ハイパーテキストualな関係を設定することもできる。

【0057】ここで既知概念と既知の事柄について更に説明する。既知概念は基本的には実質的な意味を持つ単語としての内容語のリストである。例えば富士通が何を
10 している会社なのか知っている場合には、既知概念として「富士通」を与える。そうすると要約対象文書が富士通は日本の計算機メーカーである。その富士通で今・・・が行われようとしている。・・・のように始まっている時、通常は第2文の先頭の「その富士通」に「その」という照応表現が含まれているため、第1文も要約に含むように処理される。しかしながら「富士通」は既知の概念であることが知られ、また第1文は後述するような名詞文であり、富士通の紹介（属性の定義）をしている
20 だけであると計算機の処理でも簡単に判断可能なので、本実施例では第1文は抽出されないことになる（第2文の「その」は削除される。）。

【0058】但し、このような単純な方法をとる場合には、次のような文書を対象にするとときに不都合がおきることがある。「富士通はもとと交換機をつくる会社だった、その富士通が、大型計算機で世界二位の地位をしめるようになり、今ではパソコンメーカーとして知っている
30 人の方が多いくらいだろう。だからNTTと富士通の組み合わせを不思議に思う人もいるかもしれないが、NTTと富士通の関係は浅からぬものなのである。・・・」例えば利用者がパソコンの富士通しか知らない場合、第3文を要約に含める時には、第1文も共に含めないと理解が困難となる。これを回避するためには、利用者が「富士通」についてどのような事柄を知っているのかまで指定する必要がある。例えば「富士通は日本の計算機メーカーである」、「富士通はパソコンを作っている」ことまでは知っているというように指定しなければならない。更に一步進むと、「富士通は交換機を作っていた」ことは知っているも、「富士通が今でも交換機を作っている」のか、あるいは「富士通がこれからも交換
40 機を作り続ける」のかは知らないということがあり得る。これが本実施例において事柄の既知性を取り扱うことの意味である。

【0059】このような意味で、本実施例では既知概念としては主に専門用語のリストを与える。専門用語は専門的な概念に名前をつけたものであり、分野を誤らない限りは、それが既知かどうかは容易に決められる。既知の事柄に関する知識としては、上の例のような単文の形、または単文の内容に相当するものを、フレーム表現など各種の形式で表現したものを与える。

【0060】また既知の概念は、補足的・説明的な記述が必要であるか否かの判定に用いられるために、既知概念は文分割・依存関係設定部13に与えられる。すなわちある部分を要約に含める時、それを説明しているような別の部分も含めるか否かの判定が文分割・依存関係設定部13の役割であって、ある単語が既知の概念であるかどうかは、その判定に強く影響を与えるためである。

【0061】文書要約装置の構成の説明に続いて、この装置内での本発明に特有の構成要素としての文分割・依存関係設定部13と、文選択部14の処理について説明する。図4は、文分割・依存関係設定部13による、文分割・依存関係設定処理の詳細フローチャートである。

【0062】図4においては、文分割処理と依存関係設定処理とが行われるが、文分割処理は最終的に依存関係付述語句リストを作成する文分割・依存関係設定部13の処理の前半部分である。但し図4においては、文を認定し文と文の間に依存関係を設定する処理の途中に文の内部構造を解析し、文内部の述語句間に依存関係を設定する処理が挟み込まれているので、文分割処理と依存関係設定処理とを単純に2分することはできない。

【0063】文分割処理とは、文解析部12によって単語列に変換された文書の内容を、文選択部14による文選択処理における基本単位（述語句あるいは文）に分割して、分割された基本単位を要素とするリストとしての述語句リストを作成する処理である。この文分割処理は、図4において、ステップS2による先頭の文を取り出すという処理と、点線で囲まれた述語句への分解処理の中で実行される。

【0064】図4において処理が開始されると、まずステップS1で最終的に作成されるべき述語句リストの内容がクリアされ、ステップS2で先頭の文が取り出され、ステップS3で文が取り出せたか否かが判定される。文書からの文の取り出しは、例えば見出しをそれだけで1つの文と見なすなど、文書の論理構造も考慮しながら、ピリオドなどの文の終了マークを手がかりとして行うことができる。

【0065】文が取り出せたと判定されると、ステップS4で文の構文的依存構造の解析が行われる。この構文的依存構造を求める方法としては、句構造文法によるものの、係り受け解析によるものなど様々な方法が知られているので、それらのいずれかをを用いればよい。

【0066】構文的依存構造の解析結果を基にして、ステップS5で取り出された文の述語句（単文）への分解が行われる。述語句とは、1つの述語とそれに支配される名詞（主語を含む）を基本として構成される句であり、文に含まれる単文に相当する。日本語なら用言、英語なら動詞などのような依存構造中の述語が取り出され、それに依存している要素のうち述語でないものを加えたものが述語句である。接続詞や助詞、前置詞などの機能語や、機能語に相当する表現はその前後の自立語

(内容語)とまとめておけばよい。

【0067】なお単独で名詞を修飾する形容詞などの修飾要素は、被修飾要素と一緒にまとめてもよく、独立した述語句としてもよい。但し好ましくは修飾要素の語彙的性質や、修飾要素と修飾要素の組み合わせの種類などによって、独立した述語句とするかどうかを決定する。

【0068】依存構造の解析方法としてどのような方法を用いるとしても、述語句への分解のコストはかなり高くなるが、述語句を単位として要約することにより、長い文の場合などでも簡潔な要約が生成可能となり、また高度な意味処理を行うために既知の事柄を図5に示すような格フレーム形式(フレームとは、属性名(スロット名)と属性値(フィラー)の組の列であり、知識表現法として周知のものである。)で与えたり、図6に示すような意味ネットワーク表現で与える場合などにおいて、既知の事柄情報と要約の単位との照合が簡単になるというメリットがある。

【0069】なお図6の意味ネットワークにおいて、アンダーラインは意味を表す単位となるシンボルを示し、アンダーラインのない矢印付の単語は関係を表す。図6ではシンボルが日本語で表現されているが、例えば「発表する」に対して英語の“announce”を対応するシンボルとしてあらかじめ定義しておくことにより、日本語だけでなく英語に対する情報としても使用することができる。

【0070】以上の理由から、要約の単位としては必要に応じて述語句、あるいは文のいずれかを使い分けことが望ましい。文を単位として要約を行う場合には、図4において点線で囲まれた述語句への分解処理、すなわちステップS4～S12を省略することができる。この部分には、前述の述語句への分解の部分と、次に述べる述語句間の依存関係の設定処理が含まれている。

【0071】次に図4における依存関係設定処理について説明する。依存関係設定処理は、文分割・依存関係設定部13における処理、すなわち図4の処理の後半部分である。この処理では、図3において要約プロセス制御部10から与えられる既知概念のリスト、および文書構造解析部11からの出力であり、文書構造の上から推定される構成要素間の依存関係に基づいて、概念の既知性基準による制約を、述語句リスト中の要素(述語句あるいは文)の間の依存関係の形で付与する処理が実行される。この処理は、文分割処理によって取り出された文毎に順次実行される。

【0072】構成要素間の依存関係は、他の文や句(依存先の文や句)と一緒に抽出した方が要約が読みやすくなるような文や句(従属文・句)に対して設定されるものである。そのような依存関係の設定対象には以下のようなものがある。

(1) 文に含まれる従属文

注目要素を含まない従属文を省略し、注目要素を含む従

属文を主文と一緒に要約に含めるようにする場合に、従属文から主文への依存関係を設定する。長い文の多い特許公報などを要約する場合に有効である。

(2) 前後の文に強く依存している文

逆接の接続詞(「しかし」など)を文頭に含む文などについて、注目要素を含まない場合には省略し、注目要素を含む場合には前の文も必ず一緒に要約に含めるようにする場合に、従属文(この例では後ろの文)から依存先の文(この例では前の文)への依存関係を設定する。短い文を積み重ねて書かれた論文などの場合に有効である。

(3) 見出しのついている部分に含まれる文

章節などに分かれている文書で、章節内の文に注目要素が含まれている時に章節の見出しを必ず一緒に要約に含めるようにする場合に、章節内の文から章節の見出しへの依存関係を設定する。マニュアルなどの構造化された長めの文書から、知りたい事柄がどこに書いてあるのかを調べる場合などに有効である。

(4) 主題となっている語句が既知でない文

特に動詞(の過去形「～た」)の文の主題となっている語(文頭の「～は」など)が要約中で初出の場合、その語が初めて出現した文(初出の文)も一緒に要約に含めるようにする場合に、主題が既知でない語の文から初出の文への依存関係を設定する。経済関係の雑誌などの長めの記事の場合に有効である。これは、概念の既知性基準による処理の一つである。

【0073】要約対象がかなり長い場合には、初出の文のかわりに、既知でない主題の語が、作成者注目情報とともに現れている文で近傍に出現するものを用いたり、既知でない主題の語が、必須格(「が」「を」「に」)を伴って出現している文で近傍に出現するものを用いた方がよい場合もある。なお、語の一致の判定は、同一表記の語を一致とするだけでなく、略語と正式名称との一致や同(類)義語動詞の一致などについても行うことが望ましい。英語でも、固有名詞が主語となっている動詞文などについて、同様な処置が有効なことがある。

(5) 照応表現を含んでいる文

日本語のコソアド語(「これ」「この」「こう」など)や英語の3人称代名詞などの照応表現が登場する文に注目情報が含まれていれば、先行表現が含まれている部分も一緒に要約に含めるようにする場合に、照応表現を含んでいる文から先行表現を含んでいる部分への依存関係を設定する。照応表現を含んでいる文にペナルティを与え、選択されにくくすることの方が有効な場合もある(特に、日本語の名詞文(「～は～だ」)や形容詞文(「～は～い」)など)。英語の場合や翻訳調の日本語文書の場合に有効である。ここで、照応表現を含んでいる文にペナルティを与えるということについて、名詞文・形容詞文に関する例をあげて補足説明する。「名詞文」「形容詞文」とは、形式的には文の述語が「名詞+

だ／です」、あるいは形容詞・形容動詞の文のことで、典型的には主題を表す「～は」の文節（俗にいう主語）を含むものである。内容的には、主題となっているものの性状や、主題に関する話者の判断を表している文であり、品定め文などとも呼ばれる（例えば三上章『現代語法新説』くろしお出版1972年）。本明細書で「名詞文」「形容詞文」というのは、内容的な観点からの表現で、必ずしも「名詞＋だ」や形容詞を述語としている文全てを対象としているのではない。逆に述語が動詞であっても「彼はよく働く」のように性状を表現する文であれば同様に考えることができる。

【0074】名詞文・形容詞文は、文書の中では図7に示す例（翻訳記事）のように、すぐ後で述べることを導入するために用いたり、話題のつながりを示すために使われることがよくある。このようなものは表現が抽象的であったりして、前後を見ないと具体的にどういう事柄を述べようとしているのかを掴めないものが多く見られる。そのため、そういう文だけを抜粋しても、表現の意図の理解が難しく、要約に採りあげるには意味のないことがある。例えば、図7の初めの例（アンダーライン部）であれば（典型的な形容詞とはいえないが）、後続の文を採りあげた方が、その記事が何について書いたものなのかを把握する助けとなる。次の例では、その直前の文およびそのさらに前の文を採りあげた方がよい。

【0075】本発明のねらいは、新聞や雑誌の記事やマニュアル類から必要な事柄（知識）を素早く見出すことにある。そこで、図7の例のような、話題にまとまりをつけたり、展開のつながりにするような文の価値は低くなる。特に照応表現がこのような文に現れた場合には、前後で述べている事柄にニュアンスを追加するようなものが多いので、ペナルティを与えて選択されにくくする、というのが照応表現にペナルティを与える趣旨である。

【0076】以上の（1）～（5）のような文に対応して、依存関係またはペナルティを設定することによって、要約の可読性を高めることができるが、実際にはそれ相応の計算コストが必要である。特に（4）の主題となっている語句が既知でない文、（5）の照応表現を含んでいる文の依存関係に対しては、その関係を完全に処理する技術が存在せず、不適当な依存関係が設定されることがあって可読性を損なう場合が考えられる。そこで本実施例においては、基本的な処理の流れを示すために、

（1）の従属文、（3）の見出しのついている部分に含まれる文、（4）の主題となっている語句が既知でない文について依存関係を設定し、（5）の照応表現を含んでいる文についてはペナルティを設定する場合について処理を説明する。なお、ペナルティについては、後述する文選択部14の処理において、文選択の基準となる情報量を、照応関係を含む文については通常の場合より減少させることによって、照応関係を含んでいる文が選択

されにくくするような処理が実行される。

【0077】図4のステップS5において文が述語句（単文）に分解された後に、ステップS6で構文的依存構造において依存関係にある述語の間に依存関係が設定され、ステップS7で他の述語句に依存していない述語句が主述語句として設定され、ステップS8で主述語句が述語句リストに追加される。なお前述のように文単位に処理を実行する場合には、これらの処理は省略され、単に文全体が主述語句とされる。主述語句とは、後続の処理において文と文の間の依存関係を設定する際に、依存先となるものである。

【0078】図8は述語句への分割と依存関係設定の例である。同図(a)において、文1における述語句2は述語句1に依存しているという依存関係が設定される。また文2に対しても、同様に述語句2が述語句1に依存するという依存関係が設定される。いずれの文においても述語句1が主述語句である。このように、依存関係は、構文的依存構造における述語間の関係が直接的な場合（文1）あるいは間接的な場合（文2）のいずれについても同様に設定される。すなわち、文1においては述語句2の述語「引いた」は接続助詞「ので」を介して直接的に述語句1の述語「休んだ」と関係しており、一方文2においては述語句2の述語「送ってくれた」は名詞「手紙」を介して間接的に述語句1の述語「しまった」と関係しているが、どちらも同じように依存関係が設定される。

【0079】図4のステップS8で主述語句が述語句リストに追加された後に、ステップS9～S12で文の代表句が決定される。代表句とは、その文に依存先がある場合に後述のステップS14で設定される依存関係の起点となる句のことである。文を単位として処理を行う場合には、ステップS9～S12の処理も不要であり、文全体を代表句（かつ主述語句）とすればよい。

【0080】まずステップS9で文に主題句が存在するか否かが判定される。主題句がある場合には、ステップS10で主題句が分離され、主題句と主述語句との間に依存関係が設定される。主題句とは、日本語の主題マーカー（「は」など）のついた体言句のことである。

【0081】図8(b)は主題句分離後の依存関係を示している。文1に対しては、主題句は「太郎は」であり、述語句2は述語句1に依存し、述語句1は主題句に依存するという関係が設定される。文2に対しては、主題句は「花子は」であり、同様に述語句2は述語句1に依存し、述語句1は主題句に依存するという関係が設定される。このように分解された述語句・主題句は、後述の文選択部14で依存関係に従って再構成され、要約に取り入れる場合には依存先とまとめた形で取り入れる。図8(b)の文1を例にとると、要約に取り入れられる可能性のあるのは、「太郎は、学校を休んだ」（主題句＋述語句1）あるいは「太郎は、風邪を引いたので学校を休ん

だ」(主題句+述語句2+述語句1)のいずれかである。

【0082】ここで主題句を分離する理由は、後続の文の中に主題句に依存する(主題句について述べている)文が含まれることが多く、そのような文を要約に含める時に主題句のみを含めて可読性を高めることができるからである。例えば、「太郎は学校へ出かけた、途中で犬に出会った。」の第2文「途中で犬に出会った。」を要約に含める時に、省略されている主語を補って「太郎は、…途中で犬に出会った。」(第1文の主題句+第2文の述語句)として可読性を高めることができる。この場合、主題句を分離しないと、主題句に続く主述語句も一緒に、「太郎は学校へ出かけた。途中で犬に出会った。」全体を要約に含めざるを得なくなる。なお、この処理は、照応表現処理の一貫として後述のステップS15で第2文に対応する述語句「途中で犬に出会った。」から第一文の主題句「太郎は」へ依存関係を設定することで実現される(説明中の実施例ではこの処理は行っていない)。

【0083】ステップS10で主題句の分離と依存関係の設定が行われた後に、ステップS11で主題句が文の代表句とされ、ステップS13の処理に移行する。ステップS9で文に主題句が存在しない場合には、ステップS12で主述語句が文の代表句とされた後に、ステップS13の処理に移行する。ここで代表句とは、文に含まれる主題句および述語句の中で、他の述語句に依存していないものを意味する。すなわち主題句が分離された文については、主題句が代表句であり、それ以外の文については主述語句が代表句である。

【0084】ステップS13およびS14では、文書構造解析処理で設定されている依存関係が述語句の間の関係に変換される。この処理は文が見出しなどに従属している部分(本文などの依存ブロック)に含まれている場合にのみ実行される。ステップS13で処理中の文が依存ブロック内の要素であるか否かが判定され、要素であると判定されると、ステップS14で処理中の文の代表句と文書構造上の依存先(ブロックの依存先)に対応する主述語句との間に依存関係が設定された後に、ステップS15の処理に移行する。ステップS13で依存ブロック内の要素でないと判定されると、ステップS14の処理を行うことなく、ステップS15の処理に移行する。なおここでは典型的な処理のみを記述しているため、処理中の文より後にある部分に処理中の文が依存するときに依存関係を設定するステップを含んでいない。そのような処理が必要なときには、依存先を指定する条件と従属文の代表句とを記憶しておき、その条件に一致する文を処理する時点で依存関係を設定すればよい。

【0085】図4の最後のステップ、すなわちステップS15では概念の既知性に基づいた依存関係やペナルティの設定処理が実行される。本実施例では、例えば動詞

文の主題語が既知でない場合に、主題語が初めて出現した文を依存先とする依存関係の設定処理と、照応表現を含む文にペナルティを与える処理が実行される。この処理の後に再びステップS2の処理に戻り、次の文の取り出しが行われ、ステップS3で文が取り出されたと判定されると、ステップS4以降の処理が繰り返される。ステップS3で文が取り出されなかったと判定された場合には、文分割・依存関係設定処理を終了する。

【0086】図4の文分割・依存関係設定処理が終了すると、図3の文選択部14による処理が行われる。この文選択処理は、文分割・依存関係設定部13から出力される依存関係は述語句列を対象として、要約に含めるべき重要な述語句を選択し、要約に含まれる述語句の選択結果リストを作成するものであり、その処理のフローチャートは図9に示される。

【0087】図9において、注目情報は注目概念リストとして扱われている。注目概念リストは、具体的には前述のように、注目情報として与えられた重みなしの名詞のリストである。これを注目概念リストと呼ぶのは次の理由からである。本実施例では、各文(述語句)を要約に取り入れるかどうかを判定するのに、文に含まれる注目情報の量を判定の第一の基準として用いており、注目情報量としては、注目情報として与えられた名詞が各々の文にいくつ含まれるかで計算している。注目情報量の計算においては、注目情報として与えられた名詞のリストに含まれる語と字面が一致する文中の語を数えてもよいし、あるいは、「百貨店」と「デパート」、「パソコン」と「パーソナルコンピュータ」のように同じ概念をあらわす語は同一とみなして数えてもよい。この意味で、図9では、注目語リストではなく注目概念リストという言葉を用いている。

【0088】図9において処理が開始されると、まずステップS20で注目概念リストが作成され、ステップS21で選択結果リストがクリアされる。注目概念リストは、要約プロセス制御部10および文解析部12から与えられる名詞リストに文書構造解析部で作成者注目の印がつけられた部分に含まれる名詞を加えたものである。要約プロセス制御部10から与えられる注目情報は、基本的には利用者注目情報であるが、複数の文書の比較のために要約を実行する場合には、比較対象となっている別の文書の作成者注目情報も含まれている。

【0089】つづいてステップS22で必須出力句リストが空であるか否かが判定され、空でない場合にはステップS23~24必須出力句を選択結果リストへ加える処理が実行される。ここで、必須出力句とは、利用者が要約プロセス制御部10を通じて要約に必須で含めるように指示した文書中の要素(見出しなど)と対応する述語句のことである。具体的には、文分解・依存関係設定部13から出力された述語句のうちで、文書構造解析部11によって必須出力の印がつけられた部分と対応する

もののことである。ステップS23で、先頭の必須出力句を取り出し、(その先頭の句を必須出力句リストから除いて、)取り出した句を選択結果リストに追加してから、ステップS24で、選択結果リストに追加した句の中の事柄を要約プロセス制御部から送られた既知の事柄のリストに追加し、ステップS22の処理へ戻る。

【0090】この時選択結果リストに加えた述語句の中に含まれる注目概念を注目概念リストから除くこともできる。但し、通常、必須出力要素として指定されるものは見出しなどであって、要約の核となる概念を含んでいても、完全な文の体裁をとっていないことが多いので、前述の述語句の中の注目概念を注目概念リストから除かない方がよい。述語句が完全な文の体裁をとっている場合に限って、述語句の中の注目概念を注目概念リストから除くという方法が有効である。

【0091】ステップS22～S24の処理が繰り返され、ステップS22で必須出力句リストが空になったと判定されると、ステップS25で文分割・依存関係設定部13から出力された述語句列のうちで、必須出力句以外の全ての述語句が選択候補とされて、選択候補リストが作成され、続いてステップS26で選択候補リストの中の全ての句に対して注目情報量が計算される。

【0092】この情報量の計算処理では、典型的にはそれぞれの述語句に含まれる注目概念(名詞)の数が求められる。この時注目情報量を計算しようとしている述語句に依存先があり、依存先がまだ選択結果リストに含まれていない場合には、依存先も含めて注目概念の数を数えることにする。また依存先が複数ある場合には、依存先の注目情報量を先に計算して、注目情報量が最も大きい依存先を選ぶことを原則とする。なお述語句に出現する注目概念の種類(名詞の異なり数)と、延べ数を記憶しておくことが好ましく、また注目概念に重みが与えられている場合には、その重みを乗算して名詞の数を数えるものとする。

【0093】このように依存先を含めて注目情報量を計算する理由は、前述の概念の既知性基準と関連している。すなわち概念の既知性基準に従えば、例えば文書に固有名詞が繰り返し出現する場合、2番目に出現した部分を要約に含める時には、初めて出現した部分も要約に含めるように処理することが必要となるためである。すなわち文分割・依存関係設定部13によって、2番目の文から1番目の文に対して依存関係が設定されていることから、文選択部14では2番目の文の重要度、すなわち注目情報量の計算は1番目の文とまとめて行われることになる。このような処理については具体例を用いて更に後述する。

【0094】続いてステップS27で注目情報量が0の述語句が選択候補リストから除かれた後に、ステップS28で選択候補リストに残っている全ての述語句に対して新情報量の計算が行われる。新情報量とは、利用者に

とって既知ではなく、またすでに選択された述語句には含まれていない事柄に関する情報の量である。この新情報量の計算について図10の例を用いて説明する。

【0095】図10において述語と名詞の組を新情報とし、その組の個数として新情報量を計算する。本文には7個の事柄が含まれるが、そのうち2個は見出しと共通である。見出しを必須出力要素として、図9のステップS23で選択結果リストに追加すれば、本文に含まれる新情報量は5個となる。この例では、述語と名詞の組み合わせを認定するというやや複雑な処理を必要とするが、単純に注目概念(注目語)以外の名詞の数を数えるような、より簡単な方法を用いることも可能である。

【0096】このように、新情報量の計算としては、概念の列として事柄の情報をモデル化し、すでに選択された述語句には含まれていない事柄の数を数えて新情報量とすればよい。別の方法としては、いわゆる5W1Hのような形式で事柄情報をモデル化し、図5に示したようなフレーム表現によって述語句と比較して、既知の事柄と一致しない述語句の数を新情報量とすることもできる。あるいは5W1H形式のモデルを用いた場合の新情報量を第1新情報量とし、単純な新情報量を第2新情報量として、組み合わせて用いることもできる。新情報量の計算も、注目情報量の計算と同様に、依存先があれば新情報量が最も大きな依存先を選んで、依存先を含めて計算するものとする。またペナルティが与えられている述語句に関しては、そのペナルティ分だけ新情報量を減算するものとする。

【0097】ここでペナルティについて更に説明する。照応表現などに関してペナルティを与える文としては、依存先の文であっても、依存元の文であってもよい。単に新情報量を減算するだけのことである。例えば次の2つの文についてのペナルティを説明する。

【0098】第1文：昨日¹ 新宿² で田中³ さんに会って、こんな話⁴ を聞きました。

第2文：昨日¹ 田中² さんに会って、車³ の話⁴ を聞きました。

この例では上に数字の書かれている名詞の数は同じである。しかしながら第1文の方には「こんな話」という照応表現があり、この文だけを取り出すと、どんな話であるかが不明である。照応表現に対するペナルティとは、このような場合に第2文が優先的に選択されるように、第1文に対する新情報量を、例えば名詞の個数で0.5個分減点するものである。

【0099】簡単な例として、新情報量を文に含まれる注目概念(注目語)以外の名詞の数として計算する。例えば「田中(さん)」が注目語であれば、「こんな話」に対して名詞0.5個分のペナルティを与えると、第1文の新情報量は2.5となり、第2文(新情報量3.0)の方が新情報量が多くなり、後述するように第2文の方が優先的に選択される。しかしこの例で注目語が

「田中（さん）」と「新宿」の2つであれば、第1文の方が注目情報量が多くなるため、後述するように、ペナルティにかかわらず第1文が選択されることになる。

【0100】ステップS28の処理が終了すると、ステップS29で新情報量が0の述語句が選択候補リストから除かれた後に、ステップS30で選択候補リストが空になったと判定されるまで、ステップS31～S36の処理が繰り返される。

【0101】まずステップS31で注目情報量が最大の述語句を選び、その中で最大の新情報量を持つ述語句が出力句とされ、ステップS32でその出力句が選択候補リストから除かれて、選択結果リストに追加される。この時述語句に依存先があり、その依存先がまだ選択結果リストに加えられていなければ、その依存先も一緒に選択結果リストに追加する。なお情報量が全く同じ述語句が存在する場合には、それらの述語句を全て一度に追加することを原則とするが、別の方法として最も要約対象文書の先頭に近い述語句を選ぶなど、述語句の出現位置などによって1つに選択してもよい。

【0102】このように文選択処理において事柄の既知性に関しては新情報量の差として取り扱われ、注目情報量が同じ文がある場合、新情報量の多い文が選択される。そこで注目情報量が同じ文がなければ、事柄の既知性基準は使用されないことになる。

【0103】続いて図9のステップS33で出力句、すなわち選択結果リストに追加された述語句に含まれる注目概念が注目概念リストから除かれ、その結果を用いて選択候補リストに残っている全ての述語句に対する注目情報量の再計算が行われ、ステップS34で再計算された注目情報量が0の述語句が選択候補リストから除かれる。この注目情報量の再計算は前述と同様に行ってもよいが、例えば注目概念と述語句の関係をあらかじめ記憶しておいて、リストから除かれた注目概念を含む述語句と、選択結果リストに追加された述語句に依存している述語句だけを対象に計算を行うこともできる。

【0104】ステップS34の処理に続いて、ステップS35で出力句、すなわち選択結果リストに追加された述語句に含まれる事柄の情報が既知事柄リストに追加され、選択候補リストに残っている全ての句に対して新情報量の再計算が行われる。この再計算も前述と同様に行ってもよいが、例えば事柄と述語句の関係をあらかじめ記憶しておき、既知事柄リストに加えられた事柄を含む述語句、選択結果リストに追加された述語句および注目情報量に変化があった述語句に依存する述語句だけを対象に計算を行うこともできる。

【0105】図9のステップS36で新情報量が0の述語句が選択候補リストから除かれた後に、ステップS30以降の処理が繰り返され、ステップS30で選択候補リストが空になったと判定された時点で、文選択処理を終了する。

【0106】図11は、前述のように新情報量が第1新情報量と第2新情報量とに区別されている場合の、図9のステップS31における新情報量比較の詳細フローチャートである。同図において、候補述語句Aと候補述語句Bの新情報量を比較する場合には、ステップS38でまず第1新情報量については2つの述語句のうちいずれが大きいかが判定され、Aの方がBより大きい時にはAの新情報量が大きいものと判定され、逆にBの方がAより大きい場合にはBの新情報量が大きいものと判定される。これに対して第1新情報量が同じである場合には、ステップS39で第2新情報量が比較され、第2新情報量の大きい述語句の方の新情報量が大きいものと判定され、第2新情報量が等しい場合には2つの述語句AとBの新情報量は同じと判定される。

【0107】以上においてはペナルティは照応表現に対してのみ与えられるものとしたが、前述の依存関係が設定される文(1)～(5)のうちで、(3)と(4)などについてもペナルティを与えておけば、例えば利用者に理解できない用語などの出力が抑制されることになる。その場合の処理は、照応表現についてペナルティを与える場合と全く同様に実行可能である。

【0108】あるいは図9のステップS31における新情報量の比較の後で、候補述語句の長さを含めて出力句とするか否かの評価を行うこともできる。すなわち注目情報量と新情報量とが共に等しいものについては、短い述語句を優先的に選択することにすれば、利用者に理解できない用語などの出力はある程度抑制される。更に新情報量そのものの比較に代わって、新情報量と選択述語句の長さとの比（新情報の出現密度）を用いるという方法も考えられる。

【0109】いずれにしても、このような利用者に理解できない用語などの出力の抑制の問題については、新情報量の計算とからめて処理することになる。既知概念は基本的には依存関係として取り扱われるが、事柄の既知性基準に基づく新情報量の計算にも関係している。これが、概念の既知性基準と事柄の既知性基準とを、情報の既知性としてまとめた1つの理由である。

【0110】以上において本発明における文書要約方式を、一般的に、詳細に説明したが、ここで本発明の特徴について更に説明する。本発明においては、前述のように情報の注目性基準として、利用者注目情報と作成者注目情報の2つが考慮されているために、利用者の求める情報と文書において重要な情報の双方をバランスよく取り込んだ要約を作成することができる。また情報の既知性基準として、概念の既知性基準と事柄の既知性基準とを用いることによって、簡潔で分かりやすい要約が作成されるという特徴がある。

【0111】このような前述の特徴に加えて、本実施例によればまず文書の重要性に応じて自動的に要約の長さを変えることが可能となる。従来の要約作成のアルゴリ

ズムにおいては、要約に含めるべき文の数や文字の数、あるいは元の文に対する要約の長さの比率などがパラメータとして渡される場合が多い。本実施例においては、注目情報のうちで文書の中に出現しているものの量に応じた長さの要約が作成されることになり、特別のパラメータを指定することなく、適切な長さの要約が作成される。例えば利用者注目情報を重視すれば、利用者が求める情報に関係の深い文書ほど長い要約が生成されることになる。特に長さにバラつきがある一群の文書をまとめて要約するような場合には、一般に要約の比率などのパラメータを適切に設定することが難しく、この特徴は大きな長所になる。

【0112】次に本実施例においては要約の長さに関する制約にも容易に対応できるという特徴がある。本実施例では基本的には要約の長さに関する特別な処理はなされていないが、要約の長さに対して制約が与えられた時には、その制約に容易に対処することが可能である。例えば短い要約を得たい場合には、図9の文選択処理におけるステップS30において、選択候補句リストが空にならない前に処理を途中で打ち切れればよい。これは文選択処理において、注目性基準によって重要度の高い部分から順次文が選択されていることによる。

【0113】逆に長い要約を得たい場合には、一旦図9のフローチャートに従って文選択処理を行ってから、選択されなかった部分について図9のフローチャートによる処理を繰り返すことによって、適切な長さの要約を作成することができる。これは図9のステップS31で最初の処理のフローにおいて、注目情報量が最大の句の中で新情報量が最大の句だけが選ばれているために、2回目の処理のフローにおいては次に新情報量が大きい句が選ばれて、出力句とされることによる。すなわち事柄の既知性基準によって、冗長な出力が抑制されているという本発明の特徴がある意味では逆に生かされて、適切な長さの要約の作成に寄与することになる。あるいは、選択処理を繰り返す際に、前回の選択処理で得られた要約中の全名詞を次の注目情報とするなどという方法でも、効果的に関連性の高い部分を順次取り入れて、要約を拡大することも可能である。

【0114】更に本実施例によれば、要約に関する様々な、その他の制約にも容易に対応できるという特徴がある。本実施例においては要約の満たすべき要件を、情報の注目性基準および既知性基準という2つの基準に対応して整理して用いることができ、いろいろな要求に応じた文書要約装置の動作を拡張することが容易である。例えば図3において、利用者の嗜好特性や利用者の知識のようなメモリに格納されている情報を、利用者の要求に応じて様々な面から再整理して、要約作成の制約として用いることができる。また、2つの文書の作成者注目情報を、互いに別の文書を要約するための利用者注目情報のような形で与えて、要約を作成すれば、それぞれの文

書で共通して述べられている事柄や片方にしか述べられていない事柄のうち、どちらかの文書の作成者が重点を置いているものが抽出されることになるので、そのような要約を文書の比較情報として用いることも可能である。このように要約装置の基本構成を変えことなく、各種の要求に対応することが可能である。

【0115】続いて本発明の文書要約方式を用いた要約作成の具体例について説明する。図12は本明細書の〔従来の技術〕などで参照した特許公報の抄録を対象として、次の検索式と見出し（発明の名称）を注目情報として、要約としての抄録抜粋を作成した結果を示す。

【0116】（要約OR抄録OR読解OR閲覧）AND（文書ORDキュメント）

なお最後の特開平7-445666だけに対しては、検索式は以下のものである。

【0117】（要約OR抄録OR読解OR閲覧）AND（文章ORテキスト）

図12において、要約としての抄録抜粋の中の〔 〕でくくって表示されているものは注目語である。また“特徴語”は要約対象文書の中に含まれていた注目情報である。ここで特徴語は注目語の部分集合になるが、機能的には別のものである。特に利用者が注目している語の中で、ある文書に出現したものは、その文書に対するキーワードには含まれていなくても、利用者にとっては意味が深い、その文書の特徴と考えられる。

【0118】図12における要約作成の条件は①利用者注目情報として検索式に含まれている名詞を用いる。②作成者注目情報として見出し（発明の名称）に含まれている名詞を用いる。③概念の既知性基準は用いない。④既知の事柄情報としては選択された要約部分に含まれている名詞を用いる。すなわち候補述語句中に存在し、まだ要約に含まれていない名詞の数が新情報量（異なり数が第1新情報量、延べ数が第2新情報量）とされる。⑤見出し（発明の名称）を必須出力要素とする。の5つである。

【0119】図13は概念の既知性基準を用いた場合の効果を説明する要約作成の具体例である。これは経済関係のレポートに対して、見出しを注目情報として、要約を作成した例である。図13(a)は「Hancock は」という主題句に関する概念の既知性基準を用いない場合の要約であり、(b)はこの主題句に関する概念の既知性基準を用いた場合の要約である。概念の既知性基準によって追加された部分にはアンダーラインがつけられている。

【0120】図13における要約作成の条件は①利用者注目情報は指定しない。②作成者注目情報として見出しに含まれる名詞を用いる。③概念の既知性基準は(a)では使用せず、(b)では主題句に関して用いる。④既知の事柄情報としては選択された要約に含まれている名詞を用いる。すなわち候補述語句中に存在しながら、まだ要約に含まれていない名詞の数が新情報量（異なり数、延

べ数が第1、第2新情報量)とされる。⑤見出しを必須出力要素とする。の5つである。

【0121】図13の要約作成処理について更に詳細に説明する。図13における要約対象文書は以下のものである。この文書において、下線の付いた文字で以下の記号のついた文が、図13で要約として抜粋された文である。

【0122】・図13(a)でも抽出されている文(◇)
・図13(b)で新たに追加された文(☆)

○ Apple ComputerがWindows 乗り入れの強化により再建中

◇ G.Amelioは社内の機構を改革し、Macintosh の機種を半分に減らして開発費を減らし、約3,000 人の社員をレイオフしつつ、Apple Computerを建て直している。

【0123】☆ AmelioはApple の重要な地位に外部から人を入れているが、研究開発の最高の担当者chief technology officerとして53才のEllen Hancock を任命した。これはApple 再建に最も重要な地位である。これは業界で非常に尊敬されていたD.Nagel がAT&TのBell Laboratories の所長としてApple を去るまで占めていた地位である。有意義な新製品開発の経験のある人の代わりにIBMに28年間過ごしたHancock が任命されたのは驚きを持って迎えられている。Hancock はIBMという巨大な会社で育っているの、直ちに6,000 人の血気盛んな若いエンジニアやプログラマと管理スタイルの上でぶつかり合うのではないかと見られている。また、IBMでは5年位の単位で動いているのに対し、Apple は直ちに動かなければならないから、仕事のテンポが合わないのではないかと懸念もある。Hancock は数学の修士を持ち、1996年にIBMのプログラマーとして出発し、管理能力が認められて次第に昇格し、1995年にはIBM全体の約1/3を担当したが、L.Gerstnerと意見が合わず、IBMを去ってNational SemiconductorにCOOとして迎えられた。Hancock はIBMはLotus DevelopmentのNotesを買い取るべきだと長いことIBM社内で説いていたが、それが実現したのはHancock がIBMを辞めてからである。National Semiconductorでは今年2月にAmelioがAppleに移った後、後任のchief executive officerになるつもりでいたが、board of directorsが後任にLSI LogicからのB.Hallaを任命したので、National Semiconductorを辞めた。◇しかし、Hancockはソフトウェアをよく知っており、Apple再建の成否は開発が遅れ続けているCoplandにかかっているの、Hancockは妥当な人事と見られている。また、Appleでは開発管理がいい加減で製品化が不首尾になることが多かったのを是正し、また今まで大企業のマーケットに進出できなかったのを是正できるのではないかと期待されている。

【0124】その他chief operating officerとしてTexas InstrumentsからはMarco Landi, chief administr

ative officer としてはMaxtor Corp., Advanced Micro Devices, Fairchild Semiconductorなどを経たGeorge H.Scalise、またchief financial officerとしてAutomatic Data Processing Inc.とMAI Systems にいたことのあるFred D.AndersonをAmelioは任命している。

【0125】Amelioの前任者Spindlerの時はSpindlerの無理な開発促進圧力のため、開発の中心人物が多数Appleを辞め出した。しかもその多くがMicrosoftに就職している。S.Cappsは15年間AppleにいてMacintoshその他のヒット商品を考え出したが、6ヶ月前からAppleを去る決心をし、新会社を興すためベンチャー・キャピタリストと語ったところアイデアが多過ぎると言われたため、諦めてAppleの競争相手のMicrosoftに就職し、MicrosoftのInternetのツールと新しいcomputer interfaceを開発し出している。Microsoftは今年末までにNewtonのようなhand hold computer用OSのPegasusを発表する予定だが、それを使い易くするのに協力する。またCappsと共にNewtonを開発したW.SmithもMicrosoftに移った。Gatesがもっと使い易いインターフェイスを求めているからその方針に従うが、Windowsは二人にとっては初めての経験である。Windows95は複雑で同じことをするのに五つもの異なるやり方があり、単純化するのは困難と見られている。

【0126】AmelioはApple再建に妥当な手を打っていると見られるが、結果が目に見えてくるまでは少なく共1年はかかるであろう。しかし、Macintoshの売れ行きは悪化している。今年の3月に終わる四半期のAppleの売上は1年前に比べ9.7%下がって\$2.8billionになった。Macintoshの出荷は6月に終わる四半期には20%減ると業界では推定していたが、調査会社Computer Intelligenceが1,000のパソコン小売店を調べたところ、アメリカでは4月と5月における出荷台数はそれより遙かに悪く、1年前に比べて4月は29%、5月には27%減り、売上金額は4月には31%、5月には33%減っているという。この減少の一部にはパソコン業界全体の売れ行き鈍化と、春先に欠陥のあったMacintoshを多数リコールしなければならなかったのも含まれている。しかし、Merisel Inc.のようにMacintoshは今まで通り売れているとい所もある。パソコン業界全体では4月も5月も売上は10%増え、出荷台数は3%増えている。大きな減少は企業のマーケットであり、社内のパソコンの半分から1/4までがMacintoshの3,000のオフィスを対象に調べたところ、新規のパソコンの購入が2月に33%あったのに4月には14%減っている。アメリカ最大のパソコン量販チェーンのCompu USA Inc.では、Macintoshの売上は50%も下がったが、ノートブックの売れ行きはいくつかの機種がリコールされたこともあってストップしている。小売店ではどこでもMacintoshのハードウェアもソフトウェアも余り置いていないが、それは通信販売の会社から安く買えるためである。

そういう通信販売会社の最大がMicro Warehouse であり、毎年\$1.8billionの売上有るが、その半分がMacintosh のハードウェアとソフトウェアであり、夜10時までに電話やFAX で注文すれば翌日\$3の送料で配達するという優れたサービスで有名である。この会社では1月におけるMacintosh の売上は1年前に比べ60%増えたが、5月には増減なしとなった。

【0127】Dataquest によれば、世界のマルチメディアのマーケットではMacintosh が最大で、1995年には3,950,000台(1994年には2,400,000台)、次がPackard Bell で3,000,000台(2,950,000台)、それに続いてCompaqが2,900,000台(1,200,000台)、IBMが1,600,000台(800,000台)、NECが1,500,000台(500,000台)と続く。1995年のマーケット・シェアはApple が最大で、22.9%、これに続いてPackard Bellが19.2%、Compaqが11.9%、IBMが8%、NECが4.3%、Acerが2.7%、Escom が0.7%、富士通が0.6%、Highscreenが0.6%、その他が29.1%となっている。

【0128】Apple は5月にdigital cameraやその他の画像処理用装置に内蔵されるチップ上で動作する新しいOSを発表した。これはQuick Time IC(image-capture)技術の一部であり、MotorolaのチップMPC823用のmultitasking OSであり、image-capture 用装置のAPIを含んでいる。現在digital cameraの製造会社はそれぞれ独自のASICを設計し、Adobe のPhotoshopやStorm Software社のEasyPhoto といった画像処理用ソフトウェアのインタフェースを独自に開発しなければならない。QuickTime ICを使えばdigital cameraの製造会社はそういう手間が省けるので、digital cameraの値段を下げられる。Apple はこれの開発をdigital cameraや画像処理関係の会社の大手と共同で開発してきており、すでに10社以上が支持している。これを使えばパソコンなしにdigital cameraから直接Internetに画像を送れるし、カメラの中のscriptにより撮影時間を変えたり、Photoshopのフィルターを動作させることができる。

【0129】Apple はMacintosh 互換機戦略を強化しつつあるが、そのための製品を夏から出荷する。社内でTanzaniaと呼ばれる新しいMacintosh のlogic board はscalableで安く、これによってMacintosh 互換機の製造に興味を持つ会社に呼びかける。MotorolaはすでにTanzaniaの試作を済ませ、実演をした。Tanzaniaは低位と中位機種用であり、広範なオプションがある。最高200MHzまでのPowerPC603e と604eを使用出来、3個から5個までのPCIスロットがある。PS/2キーボードかADBコネクター、またEnhanced IDE (Integrated Drive Electronic)かSCSI internal hard driveの選択がある。またLocalTalk, GeoPort, SCSIなどのコネクタの他、Apple としては初めてのATADI (AT attachment packet interface) CD-ROM driveがある。Tan

zaniaはまたMacintosh の自動ejectingのfloppy driveの他、Intel のチップ使用のパソコンでは標準の手動ejectingのものもある。DIMMスロットが2個とSIMMのスロットが二つあり、最大160Mbytes のRAMが使えるEDO DRAMを使用する。互換機製造会社はTanzaniaを使用したパソコンを来年始めには出荷できる。Apple はMacintoshのライセンス戦略をMacintosh そのものの互換機からPPCPへの移行を3段階に分けて推進している。第1段階はMacintosh そのものの互換機だけであって1995年から1996年にかけてDayStar Digital Inc., Power Computing Corp., Umax Computer Corp.がPower Macintosh7500 と9500の互換機を実現した。Umax Computer Corp. は台湾のUmax Data Systems がRadius Inc. からMacintosh 互換機部門を買い取って今年1月に生まれた会社で、その最初の互換機SuperMacS900は6月始めから出荷されたが、非常な人気で生産が間に合わず、1カ月以内に\$10millionの受注がこなせないでいる。第二段階はPower Macintosh5400 とTanzaniaボードに基づいて今年夏から来年半ばにかけて実現する。この二つは共にLow End Reference Platform (最近MacOS Licensing Design, 略してMLDという)に基づいている。どちらも業界で標準になっている広範囲の論理回路や周辺機器が使えるようになっており、PPCP (以前はCommon HardwareReference Platform、略してCHRPと呼ばれていた)の狙いに近づく。PPCPではMacintosh のOSの他、OS/2, Windows 3.1, UNIX, Solarisなど広範囲のOSが使えることになっている。第三段階はPPCPに1997年半ばから1998年にかけて完全に移行する。

【0130】Microkernel に基づくCopland が来年半ばに延びたので、今までApple はCopland が出るまでSystem7.5.3 がSystemの改良の最後だと言っていたのを変更して、つなぎとしてCopland の新しい機能の一部を取り入れたHarmony と社内で呼んでいるOSを年末に発表することになった。Harmony にはInternetへのサポート、OpenDoc, Cyberdog, QuickTime2.5, QuickDraw3Dなどのグラフィックス技術、またCopland に予定されていたインターフェイスの改良などがある。またラベルの付いたフォルダーによってファイル多数の検索や管理もできる。Lockheed Martin Missiles and Space 社はMacintosh を9,500台持っているが、InternetへのサポートやOpenDoc が完成しているのなら、来年まで待たずに入手できるのは有難いと大歓迎である。

【0131】ユーザは今までのソフトウェアを変更せずHramony を使えるが、Copland の場合はソフトウェア会社は今までのソフトウェアを変更しなければならない。Copland は今はSystem8 と呼ばれている。

【0132】また今年の夏にはSystem7.5.3 のバグを直し、Duo23005やPowerBook に対する性能を改善し、Busterと社内で呼ばれて開発されてきたものを発表する。Co

puter Intelligence InfoCorp.の最近の調査では、昨年Macintoshを買った人の中で次もMacintoshを買おうと答えたのが87%もあり、パソコンの満足度では最高であった。次いでDell Computerが74%、Hewlett-Packardが72%、Acerが68%、Gateway2000が61%であった。Macintoshに満足している人々はMacintoshのOSが好きなためであるのに対し、Intelチップ使用のパソコンを使用している人々はたとえMacintoshのOSの方が好きであってもソフトウェアの互換性の立場からMacintoshに変えられないとこの調査会社は説明している。

【0133】Amelioの前任者SpindlerはMicrosoftを徹底的に敵視してGatesに会うことはなかったが、AmelioはGatesを訪問して協力を要請するという現実的な行動をしている。マルチメディアに関する標準と製品をAppleとMicrosoftの二社で共通に使用しようとしている。二社の交渉がまとまればAppleのQuickTimeの開発環境がWindows95やDirectX APIのサポートも含めNTでも使えることになり、またAppleのQuickTime Internet ExplorerがMicrosoftのInternet Explorerでも使える。◇交渉がまとまればMicrosoftはQuickTimeをInternet Explorerに組み込むようにし、同時にAppleはWindowsのマルチメディア技術のサポートを強化することになる。すでにAppleはQuickTimeがWindows使用の環境でも使用出来るようにし、DirectXのAPIの多くがQuickTimeでも使えるようになっている。今まではWindows用のQuickTimeは再生しかできなかった。Web上のビデオの60%がOpenDocで作られ、30%がMPEGにより作られるようになりつつあり、MPEGはQuickTimeでも読めるという事実の前に、Microsoftは現実的になってきた。

【0134】二社の関係は他の点でもよくなってきている。AppleはMicrosoftのBackOfficeをヨーロッパではAdvanced Workgroup Solutionsのサーバーにバンドルして売り出しており、これが成功すればアメリカでもそうすると言う。MicrosoftもOffice97の次のversionをMacintosh用にも開発すると約束するようにまで二社の関係は改善された。パソコンの需要が落ち込みつつある現在は特に他社のソフトウェアやハードウェアへの相互乗り入れは互いに利益がある。

【0135】AppleはQuickTime Internet技術をMicrosoftにライセンスを与えるばかりでなく、QuickTime VR (virtual reality) も含めることになる。これに対抗するMicrosoftのActiveMovieの技術は非常に遅れており、開発キットをソフトウェア会社に配布さえしていない。MicrosoftはActiveMovieをInternet Explorerに取り入れるのを諦めるのではないかとみられている。しかしAppleはDirectXのAPIなどMicrosoftのメディア技術のサポートを改善するためにQuickTimeを貫き直さなければならない。DirectXのAPIの大部分をAppl

eはサポートしつつある。これらAPIのうちDirect3Dは、AppleのQuickDraw3Dと真っ向から対抗する。他方IntelもInternet上のビデオ技術についてMicrosoftに働きかけている。Appleのビデオ技術よりIntelのビデオ技術の改良されたものの方がよいとMicrosoftに説いている。Appleの場合はビデオを再生し始める前に十分なデータをダウンロードしなければならないのに対し、Intelの改良した技術では圧縮技術の改良によって早くビデオを再生し始められる。GatesはInternetやintranetの将来性を見誤ったのに気がつき、今は出遅れたInternet Explorerを一刻も早く強力なものにしたいとあせっている。

(要約対象文書終わり) この文書を対象とする要約の作成における注目情報としては、見出しに含まれている名詞を用いる。すなわち注目語は“Apple Computer”、“Windows”、“強化”、および“再建”の4つである。注目情報量についても、第1注目情報量と第2注目情報量とに区分して取り扱う。第1注目情報量は文に含まれる注目語の異なり数であり、第2注目情報量は文に含まれる注目語の延べ数である。第1注目情報量、および第2注目情報量の取り扱いは、図11の新情報量の比較と同様に行われる。

【0136】説明の単純化のために、新情報量としては注目語以外の名詞(内容語)の数を用いることとし、第1新情報量は文に含まれる注目語以外の名詞内容語の異なり数、第2新情報量は文に含まれる注目語以外の名詞内容語の延べ数とする。

【0137】図14は注目情報量と新情報量の計算結果である。上の要約対象文書の先頭から文に番号をつけ、注目情報量、新情報量が0とならない文についての情報量を示す。文に対して注目語は太字の〔 〕、新情報としての注目語以外の名詞内容語は細字の〔 〕で囲んで表示されている。情報量の少数点より前の部分が異なり数であり、少数点より後の部分が延べ数である。例えば文11においては〔Hancock〕が2回出現するので、新情報量の異なり数(第1新情報量)は8、延べ数(第2新情報量)は9となる。

【0138】図14の計算結果に基づいて、図9のステップS31で文72が選択され、ステップS32で文72が選択候補リストから除かれ、選択結果リストに追加される。そしてステップS33で“Windows”と“強化”とが注目語リストから除かれ、注目情報量の再計算が行われる。

【0139】図15は再計算後の情報量を示す。この図に対しては、図9のステップS31で文1が選択され、ステップS32で選択結果リストに追加され、ステップS33で“Apple Computer”が注目語リストから除かれ、注目情報量の再計算が行われる。ここでは文1以外に“Apple Computer”を含む文は存在せず、他の文の情報量は変化しない。

【0140】続いて再びステップS31で文11が選択され、ステップS33で“再建”が注目語リストから除かれ、注目語リストが空になる。従って、注目情報量を再計算すると、選択候補リストに残っている述語句に対する注目情報量は全て0になり、ステップS34で選択候補リストの内容は空となり、文選択処理が終了する。この処理によって得られた結果が図13(a)である。

【0141】次に図13(b)の結果を得る処理について説明する。ここでは図13(a)に対する処理に加えて、以下の処理が実行される。まず第1に、主題句の中に既に固有名詞が表れた文に対しては、その固有名詞が要約対象文書中で初めて出現した文を依存先とする依存関係を設定する。但し固有名詞の場合、最初は正式名称(ここでは“Ellen Hancock”, “G.Amelio”)が用いられても、その後は省略形(ここでは“Hancock”, “Amelio”)が用いられることが多いので、そのような正式名称と省略形は同じとみなす。第2に主題句に指示詞(例えば「これ」)が出現した時には、その直前の文を依存先として依存関係を設定する。第3に依存先の文が、更に第1および第2の処理における依存関係に該当する場合には、その先の依存先の文についても同様に依存関係を設定する。

【0142】まず注目語が出現する文に対して第1〜第3の処理を行って、依存関係を設定する。図16はこの依存関係を示す。例えば文11の主題句としての「Hancockは」に関して文11から文2に対する依存関係が設定され、文2の「Amelio」に関して文2から文1への依存関係が設定されている。なお文41や文72における“Apple”や“Microsoft”も固有名詞であるが、有名な企業名であり、利用者既知概念として与えられたものとして、以下の説明を行う。

【0143】このような依存関係を考慮した情報量の計算について、文11を例にとりて説明する。図16(c)の依存関係に対応して、文11の情報量は文2と文1の情報量を含めて計算する。ここで“Apple”は“Apple Computer”と同一の語とみなす。

【0144】図17はその計算結果を示す。文11の注目情報量は文11に“再建”、文1に“Apple Computer”が含まれているため、異なり数としての第1注目情報量は2、延べ数としての第2注目情報量は“Apple”の分も含めて4になる。また新情報量は“Apple”を除き、“Hancock”, “Ellen Hancock”, “G.Amelio”, “Amelio”、および“開発”の重複を除いて、異なり数で24、延べ数で27となる。

【0145】図17の計算結果を用いて、図9のステップS31で文11が選択され、S32で文1および文2と共に選択候補リストから除かれ、選択結果リストに追加される。これによって、例えば文3の依存先は選択済となるので、次に文3が選択される場合には、その情報量は文3だけに対して計算される。その後ステップS3

3で“Apple Computer”と“再建”が注目語リストから除かれ、注目情報量の計算が行われ、文3と文2の注目情報量が0となる。この結果が図18である。

【0146】更に図18の結果によって、ステップS31で文72が選択され、選択結果リストに追加され、ステップS33で“Windows”と“強化”とが注目語リストから除かれ、注目語リストが空になるために、文選択処理が終了する。これによって図13(b)の結果が得られる。

【0147】最後に本発明の文選択方式の第2の実施例について説明する。図19は、この文選択方式としての文抽出のアルゴリズムを示す。このアルゴリズムは、新聞記事やレポートなどの見出しに含まれる名詞キーワードを用いて、その名詞キーワードを含む文を抽出して、記事などのダイジェスト情報を作成するアルゴリズムである。

【0148】図20は、図19のアルゴリズムにおける用語と、図9の文選択処理のフローチャートにおける用語との対応の説明図である。なお図9に対しては、図14に関する説明と同様に、注目情報量が第1情報量と第2情報量とに区別されている。

【0149】例えば図12と比較すると、図12においては検索式(質問文)が用いられているのに対し、図19では見出しだけが用いられる点が異なっている。見出しと質問文の違いとしては、まず第1に見出しは必須出力要素であり、見出しと全く同じ語しか含まない文、すなわち新情報量が0の文は抽出されないことと、第2に質問文(検索式)は単なる注目語のリストであり、質問文と全く同じ語しか含まない文であっても抽出される、すなわち質問文自体は選択結果リストに含まれず、新情報量が0にならないということがある。

【0150】図20において、新情報量と見出しキーワードに一致しない名詞の延べ数とを対応させることは、注目概念(見出しキーワード)と同じ文に出てくる名詞の組の数によって新しい事柄の量を求めるという考え方と同じである。すなわち(3)の比較、注目概念(正確には選択結果リストにまだ含まれていない注目概念)の数が全く同一の文について行われるため、まだ出現していない注目概念とそれ以外の名詞の組み合わせを数えることと同じになる。

【0151】図19においては見出しだけを注目概念のソースとして用いるために、注目概念同士の組は見出しにすでに出現していることになり、図9では必須出力要素としてすでに選択結果リストに含まれていることになる。それ以外の注目概念に関する名詞の組の数は、文に含まれる見出しキーワードの数と見出しキーワード以外の名詞の数との積で与えられ、見出しキーワードの数が同じとすれば、見出しキーワード以外の名詞の数だけを(3)において比較することにより、名詞の組の数の比較を行うことと同等となる。

【0152】

【発明の効果】以上詳細に説明したように、本発明の文書要約装置を用いることにより、前述の様々な特徴によって各種の効果が生ずるが、最も大きな第1の効果として、文書の有用性（関連度）の判定が容易になるという効果がある。すなわち本発明の方法によれば、利用者が注目している情報と、作成者が重点的に記述しようとしている情報との両方を要約の中に抽出することができるため、利用者が注目している情報が文書の中でどのように扱われているかが、要約を読むだけで容易に把握可能となる。すなわち、文書が利用者の目的とどの位関係があるかということを、要約から容易に判定できるようになる。

【0153】第2の大きな効果としては、要約の可読性が向上するという効果がある。概念の既知性基準によって利用者が知らない用語については、例えば追加説明と共に出力され、また事柄の既知性基準によって冗長な出力が抑制されるので、簡潔かつ分かりやすい要約が作成される。また利用者注目情報に基づいて、利用者の注意の方向を考慮することにより、利用者に不必要な情報が要約中に含まれる率を減少させることができることも、利用者にとってより分かりやすくする上で大きな効果がある。

【図面の簡単な説明】

【図1】本発明の第1の原理の説明図である。

【図2】本発明の第2の原理の説明図である。

【図3】本発明の文書要約装置の構成を示すブロック図である。

【図4】文分割・依存関係設定処理の詳細フローチャートである。

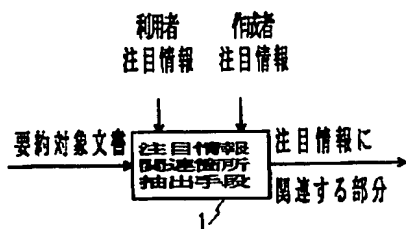
【図5】格フレームによる事柄情報の表現の例を示す図である。

【図6】意味ネットワークによる事柄情報の表現の例を示す図である。

【図7】照応表現にペナルティを与える理由を説明するための例を示す図である。

【図1】

本発明の第1の原理の説明図



【図8】述語句への分割と依存関係の例を示す図である。

【図9】文選択処理の詳細フローチャートである。

【図10】新情報量の計算方法を説明する図である。

【図11】第1と第2の新情報量を区別する場合の新情報量の比較処理フローチャートである。

【図12】特許抄録の要約結果の例を示す図である。

【図13】主題句に関する概念の既知性基準の使用の効果を説明する要約例を示す図である。

【図14】図13(a)を得るための情報量の初回計算結果を示す図である。

【図15】図14において文72が選択された後の情報量を示す図である。

【図16】図13(b)の結果を得るための依存関係の設定を説明する図である。

【図17】図16の依存関係を考慮した情報量の計算結果を示す図である。

【図18】図17において文11が選択された後の情報量を示す図である。

【図19】文選択方式の他の実施例としてのダイジェスト情報抽出のアルゴリズムを示す図である。

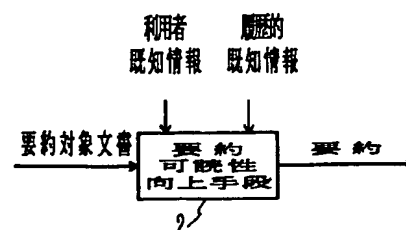
【図20】図19のアルゴリズムと図9のフローチャートにおける用語の対応を示す図である。

【符号の説明】

- 1 注目情報関連箇所抽出手段
- 2 要約可読性向上手段
- 10 要約プロセス制御部
- 11 文書構造解析部
- 12 文解析部
- 13 文分割・依存関係設定部
- 14 文選択部
- 15 要約整形部
- 16 利用者の嗜好特性
- 17 利用者の知識
- 18 閲覧履歴
- 19 入力文書（群）

【図2】

本発明の第2の原理の説明図

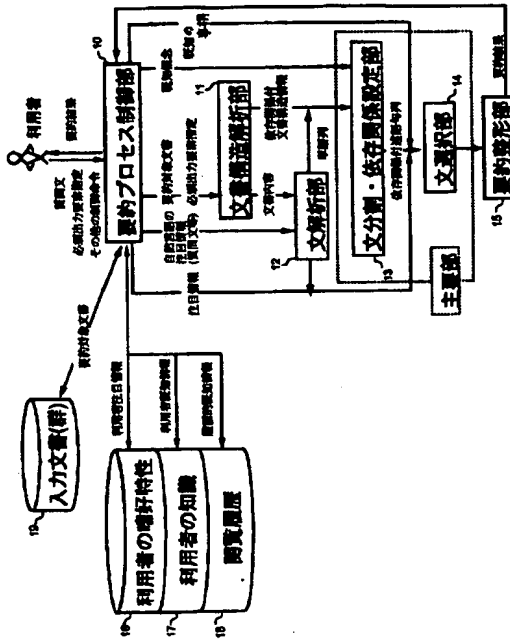


【図3】

【図 5】

格フレームによる亭柄情報の表現の例を示す図

本発明の文書要約装置の構成を示すブロック図



【図7】

動詞: 発表
 ガ格: 富士通
 ヲ格: パソコンの新製品
 時: 平成3年3月3日
 場所: 日本

照応表現にペナルティを与える
理由を説明するための例を示す図

「...この需要のトリゾは、純粋の需要ではない。顧客層となったほとんどのクライアントで、いまはコンピュータはDOS/386、ミッドレベルでUNIX、そしてチャートメントグラフィックスレベルでENTが使われている。

——日本におけるコンピュータ市場の調査は、これまで一タークに亘らぬで、だが、コンピュータのマルチメディア化の勢、テレビとの融合が進んで、日本企業特有な利点からいっても、この分野では、日本のクライアントにとっては極めて、Compaq computerは日本の市場のチャートグラフについで、日本を第2位とした。

(…は環境現象以外で言明した文がある部分を示す。)

【図6】

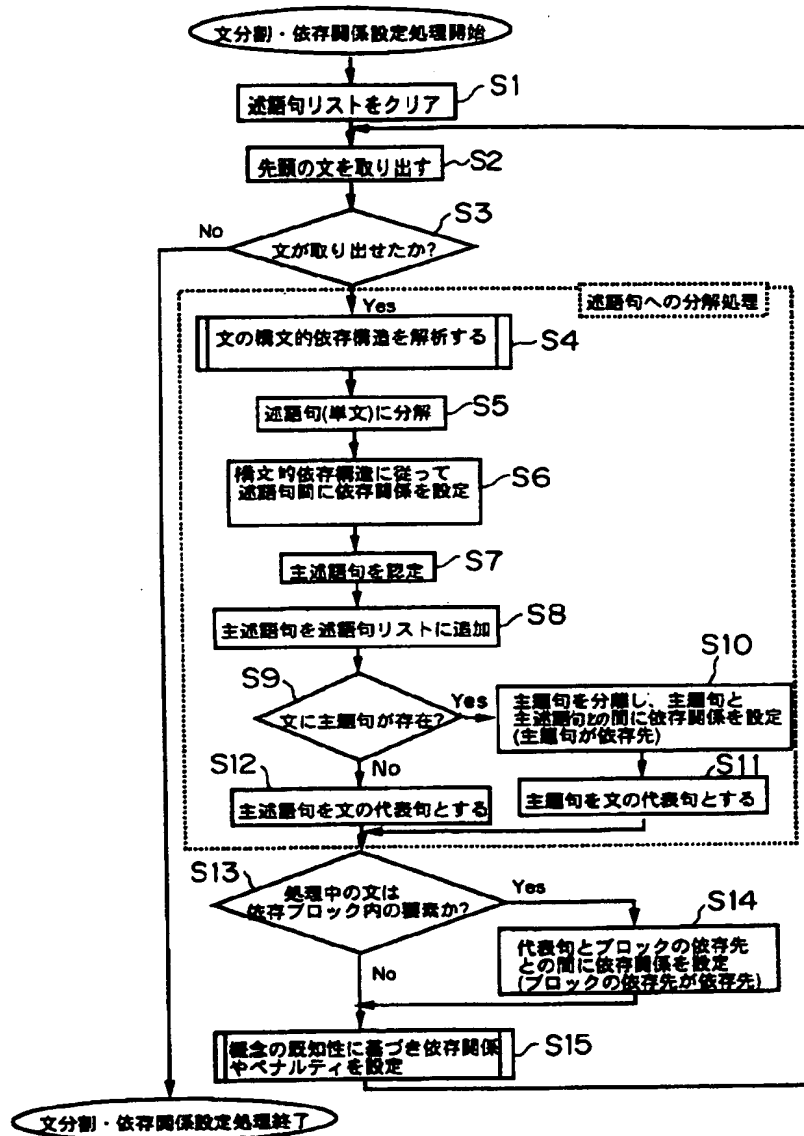
意味ネットワークによる
事柄情報の表現の例を示す図

```

発表する — 行為者 → 富士通
      |
      | 対象 → 製品 — 属性 → 新しい
      |
      | 限定 → パソコン
  
```

【図4】

文分割・依存関係設定処理の詳細フローチャート



【図8】

述語句への分割と依存関係の例を示す図

文1: 太郎は、風邪を引いたので学校を休んだ。
 述語句1: 太郎は、/ 学校を / 休んだ。
 述語句2: 風邪を / 引いたので
 依存関係: <述語句1> ←依存- <述語句2>

文2: 花子は、太郎が送ってくれた手紙を大事そうにしまった。
 述語句1: 花子は、/ 手紙を / 大事そうに / しまった。
 述語句2: 太郎が / 送ってくれた
 依存関係: <述語句1> ←依存- <述語句2>
 (a) 主題句分離前

文1: 太郎は、風邪を引いたので学校を休んだ。
 主題句(代表): 太郎は、
 述語句1(主): 学校を / 休んだ。
 述語句2: 風邪を / 引いたので
 依存関係: <主題句> ←依存- <述語句1> ←依存- <述語句2>

文2: 花子は、太郎が送ってくれた手紙を大事そうにしまった。
 主題句(代表): 花子は、
 述語句1(主): 手紙を / 大事そうに / しまった。
 述語句2: 太郎が / 送ってくれた
 依存関係: <主題句> ←依存- <述語句1> ←依存- <述語句2>
 (b) 主題句分離後

【図13】

主題句に関する概念の既知性基準の
 使用の効果を説明する要約例を示す図

Title: Apple Computer が Windows 乗り入れの強化により再建中

G.Amelio は社内の職務を改革し、Macintosh の開発を十分に知らしめて開発者を減らし、約 3,000 人の社員をレイオフしつつ、Apple Computer を建て直している。…しかし Hancock はソフトウェアをよく知っており、Apple 両派の両者は開発が遅れ続けている Copland にかかっている。Hancock は優秀な人事と見られている。…交渉がまとまれば Microsoft は QuickTime を Internet Explorer に組み込むようになり、同時に Apple は Windows のマルチメディア技術のサポートを強化することになる。

(a) 概念の既知基準を用いない場合

Title: Apple Computer が Windows 乗り入れの強化により再建中

G.Amelio は社内の職務を改革し、Macintosh の開発を十分に知らしめて開発者を減らし、約 3,000 人の社員をレイオフしつつ、Apple Computer を建て直している。Amelio は Apple の最高な地位に外務から人を入れているが、研究開発の最高責任者 chief technology officer として 63 歳の Ellen Hancock を任命した。…しかし Hancock はソフトウェアをよく知っており、Apple 両派の両者は開発が遅れ続けている Copland にかかっている。…交渉がまとまれば Microsoft は QuickTime を Internet Explorer に組み込むようになり、同時に Apple は Windows のマルチメディア技術のサポートを強化することになる。

(b) 主題句に関して概念の既知基準を用いた場合

【図10】

新情報量の計算方法を説明する図

見出し: 富士通新情報機器発売へ

本文: 富士通は三十日、新情報機器の発売を三月に開始すると発表した。

見出しに含まれる事項	本文に含まれる事項
(富士通、発表) (新情報機器、発表)	(富士通、発表) (新情報機器、発表) (発表、開始) (三月、開始) (富士通、発表) (三十日、発表) (開始、発表)

【図14】

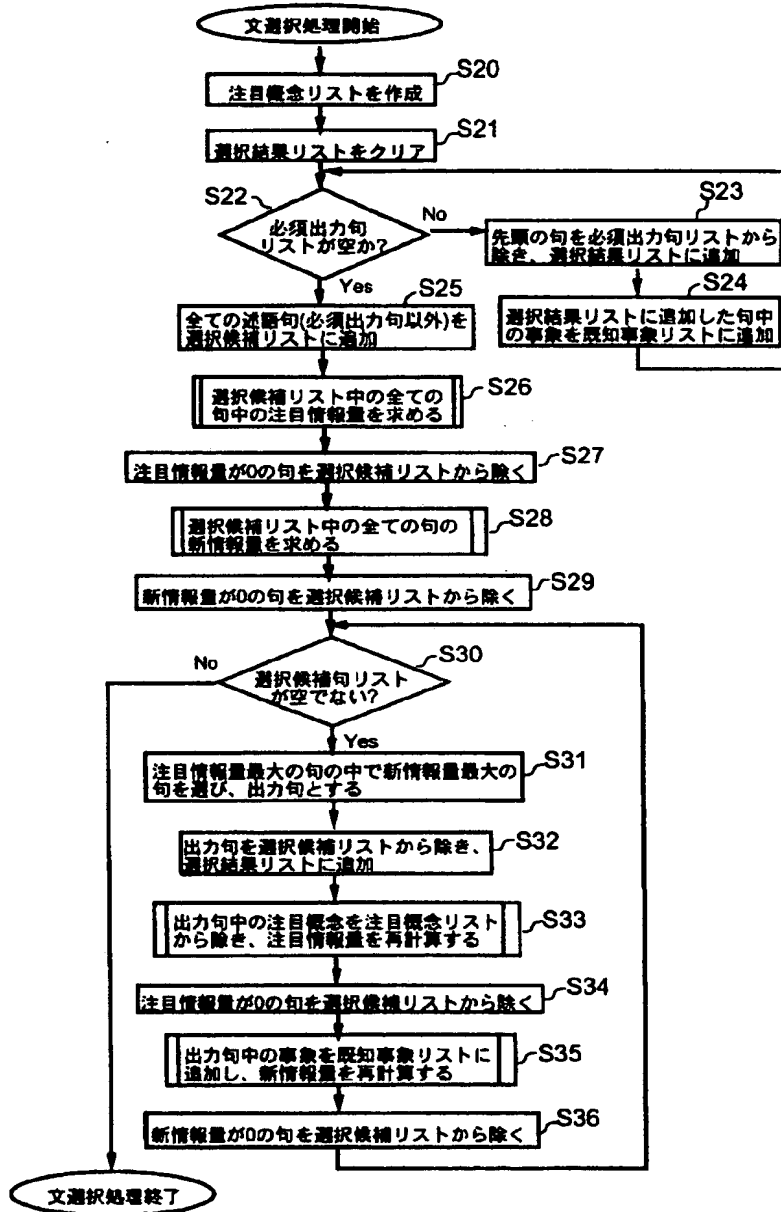
図13(a)を得るための

情報量の初回計算結果を示す図

文番号	注目情報量	文
1	1.1	(G. Amelio) は (社内) の (職務) を (改革) し、(Macintosh) の (開発) を十分に知らしめて (開発) 者を減らし、約 3,000 人の (社員) を (レイオフ) しつつ、(Apple Computer) を (建て直) している。
3	1.1	これは (Apple) (両派) に (優秀) な (人事) と (見) られている。
11	1.1	しかし (Hancock) は (ソフトウェア) を よく知っており、(Apple) (両派) の (両派) は (開発) が (遅) れ続けている (Copland) にかかっているの、(Hancock) は (優秀) な (人事) と (見) られている。
21	1.1	(Amelio) は (Apple) (両派) に (優秀) な (手) を (打) っていると思われ、結果が (目) に見えてくるまでは少なく (共) 1 年 (は) かかるであろう。
41	1.1	(Apple) は (Macintosh) (両派) 派 (両派) を (強化) しつつあるが、そのための (製品) を (開発) する。
72	2.2	(交渉) が (まと) まれば (Macintosh) は (QuickTime) を (Internet Explorer) に (組み) 込むようになり、同時に (Apple) は (Windows) の (マルチメディア) (技術) の (サポート) を (強化) することになる。
74	1.1	今までは (Windows) 用の (QuickTime) は (両派) しかできなかった。

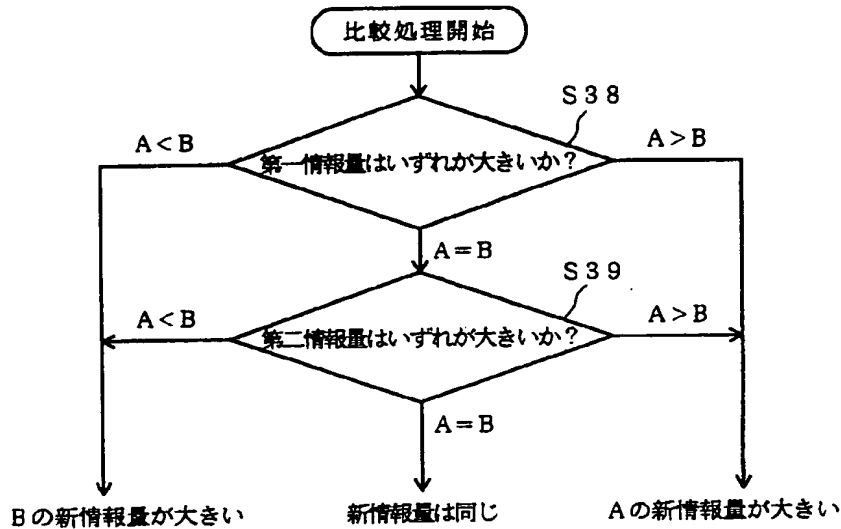
【図9】

文選択処理の詳細フローチャート



【図11】

第1と第2の新情報量を区別する場合
の新情報量の比較処理フローチャート



【図19】

【図20】

文選択方式の他の実施例としての
ダイジェスト情報抽出のアルゴリズムを示す図

```

未出現キーワード列 ← 見出しキーワード。
抽出結果 ← 空。
while 未出現キーワード列が空でない
do (重要文選択)
  未出現キーワード列に基づき文の重要
  度を評価。
  重要度最大の文 (複数可) を選択。
  (1つも選べない場合は処理終了。)
  選択された文を抽出結果に加える。
  選択された文中の見出しキーワードを
  未出現キーワード列から除く。
od
  
```

図19のアルゴリズムと図9のフローチャートにおける用語の対応を示す図

図19	図9
未出現キーワード列	注目全文リスト
抽出結果	選択結果リスト
重要度	注目情報量と新情報量
(1) 見出しキーワードと一致する名詞の異なり数	(第一注目情報量)
(2) 見出しキーワードと一致する名詞の延べ数	(第二注目情報量)
(3) 見出しキーワードと一致しない名詞の延べ数	(新情報量)

【図12】

特許抄録の要約結果の例を示す図

公開番号	特開平 06-259424
発明の名称	文書表示装置及び文書要約装置並びにデジタル複写装置
発明者	亀田 肇之(リコー)
抄録技術	…本発明による【文書】【表示】【装置】は、見出し部分と本文とから成る【文書】に対する【表示】機能を有するもので、解析手段1は見出し部分を解析し、認識手段2は、前記解析手段1により解析された見出し部分中の単語を本文中から認識する。…
特徴語	装置、表示、文書
公開番号	特開平 07-036886
発明の名称	文書を要約する方法および装置
発明者	エム マーガレット ウィズゴット/ダグラス アール カツティング(ゼロックス)
抄録技術	…本発明の【方法】は、【文書】の領域を選択することによって【要約】を自動的に作成する。…
特徴語	文書、方法、要約
公開番号	特開平 08-297877
発明の名称	主題の要約を生成する自動的な方法
発明者	フランシーヌ アール チェン(ゼロックス)
抄録技術	【課題】機械で読み取り可能なドキュメントの【主題】の【要約】を自動的に【生成】する【方法】を提供する。…
特徴語	主題、生成、方法、要約
公開番号	特開平 06-215049
発明の名称	文書要約装置
発明者	乾 隆夫/芥子 育雄/石松 隆一郎(シャープ)
抄録技術	…文書処理部4は、各文脈ベクトル間距離を参照して、【文書】に最も近い段落と【文書】に近い複数文との2種類の要旨及び【文書】に最も近い各段落毎の文と各段落に最も近い文との2種類の【要約】を生成する。…
特徴語	文書、要約
公開番号	特開平 07-182373
発明の名称	文書情報検索装置及び文書検索結果表示方法
発明者	住田 一男/三池 誠司/小野 順司/竹林 洋一/武田 公人/伊藤 悦雄(東芝)
抄録技術	【検索】結果を原文で【表示】するのではなく、利用者の所望する観点での【要約】文を提示することにより、利用者が【検索】した【文書】の内容をたやすく理解し、要不要をすばやく判定することを可能にする【文書】【情報】【検索】【装置】を提供する。…
特徴語	検索、情報、装置、表示、文書、要約
公開番号	特開平 07-044568
発明の名称	抄録作成装置
発明者	小野 順司/住田 一男/三池 誠司(東芝)
抄録技術	(J)【抄録】文中の隠蔽表現の隠蔽先或いは省略箇所表示を【抄録】文中に含ませることにより、【抄録】文の読解性、自然性の向上を図る。書式解析部1は電子化された入力【文章】を解析し、文の切れ目、段落の変わり目、章や節の構造等を解析する。…
特徴語	抄録、文章

【図15】

図14において文72が
選択された後の情報量を示す図

文 番号	注目 情報量	新 情報量	文
1	1.1	8.9	(G. Amelio)は(社内)の(報酬)を(改定)し、(Macintosh)の(報酬)を平均に属らして(報酬)量を増らし、約3,000人の(社員)を(レイオフ)しつつ、【Apple Computer】を建て替えている。
3	1.1	3.3	これは【Apple】(再読)に最も(重要)な(地位)である。
11	1.1	8.9	しかし(Hamock)は(ソフトウェア)をよく知っており、(Apple)【再読】の(経営)は(報酬)が掛け回している(Copland)にかかっている、(Hamock)は(要請)な(人事)と見られている。
21	1.1	5.5	(Amelio)は(Apple)【再読】に(要請)な(手)を打っていると思われるが、結果が(目)に見えてくるまでは少なく、1年はかかるであろう。
41	0.0	6.6	(Apple)は(Macintosh)【再読】の(製品)を(更新)しつつあるが、そのための(製品)を(更新)する。
72	...	(読みずみ)	
74	0.0	2.2	今までは【Windows】用の(QuickTime)は(再生)しかできなかった。

【図16】

図13(b)の結果を得るための
依存関係の設定を説明する図

文 番号	注目 情報量	新 情報量	文
1	1.1	8.9	(G. Amelio)は(社内)の(報酬)を(改定)し、(Macintosh)の(報酬)を平均に属らして(報酬)量を増らし、約3,000人の(社員)を(レイオフ)しつつ、【Apple Computer】を建て替えている。
2	0.0	11.1	(Amelio)は(Apple)の(重要)な(地位)に(外部)から(人)を入れているが、【新再読】の(経営)の(担当者)(chief technology officer)として(経営)の(Bill Hamock)を(任命)した。
3	1.1	3.3	これは【Apple】(再読)に最も(重要)な(地位)である。
11	1.1	8.9	しかし(Hamock)は(ソフトウェア)をよく知っており、(Apple)【再読】の(経営)は(報酬)が掛け回している(Copland)にかかっている、(Hamock)は(要請)な(人事)と見られている。
21	1.1	5.5	(Amelio)は(Apple)【再読】に(要請)な(手)を打っていると思われるが、結果が(目)に見えてくるまでは少なく、1年はかかるであろう。

(a) 依存関係にある部分
 文1 ← 依存 (Amelio)は → 文2 ← 依存 (これは) ← 文3
 (b) 文3に属する依存関係
 文1 ← 依存 (Amelio)は → 文2 ← 依存 (Hamock)は → 文11
 (c) 文11に属する依存関係
 文1 ← 依存 (Amelio)は → 文21
 (d) 文21に属する依存関係

【図17】

【図18】

図16の依存関係と考慮した
情報量の計算結果を示す図

文 番号	注目 情報量	新 情報量	文
1	1.1	9.9	(G. Amelio)は(社内)の(情報)を(提供)し、(Macintosh)の(機能)を十分に備わって(開発)費を削減し、約5,000人の(社員)を(レイオフ)しつつ、(Apple Computer)を(強化)している。
3<(21)>	2.4	18.21	これは(Applc)【所属】に最も(重要)な(施設)である。
11<(21)>	2.4	24.27	しかし(Hamrick)は(ソフトウェア)をよく知っており、(Apple)【所属】の(経営)は(関係)が(維持)されている(Copland)にかかっているため、(Hamrick)は(担当)な(人事)と見られている。
21<(1)>	2.3	12.13	(Amelio)は(Applc)【所属】に(担当)な(手)を行っていると思われるが、結果が(目)に見えてくるまでは少なく(来1年)は(かかる)であろう。
41	1.1	6.6	(Apple)は(Macintosh)【互換】な(製品)を(強化)しつつあるが、そのための(製品)を(出資)する。
72	2.2	6.8	(文部)が(とまれば(Macintosh)は(QuickTime)を(Internet Explorer)に組み込むようにし、同時に(Applc)は(Windows)の(マルチメディア)【技術】の(サポート)を(強化)することになる。
74	1.1	2.2	今までは(Windows)用の(QuickTime)は(再生)しかできなかった。

図17において文11が選択された後の情報量を示す図

文 番号	注目 情報量	新 情報量	文
1			(選択済み)
8	0.0	3.3	これは(Applc)【所属】に最も(重要)な(施設)である。
11<(21)>			(選択済み)
21	0.0	5.5	(Amelio)は(Applc)【所属】に(担当)な(手)を行っていると思われるが、結果が(目)に見えてくるまでは少なく(来1年)は(かかる)であろう。
41	1.1	6.6	(Apple)は(Macintosh)【互換】な(製品)を(強化)しつつあるが、そのための(製品)を(出資)する。
72	2.2	6.8	(文部)が(とまれば(Macintosh)は(QuickTime)を(Internet Explorer)に組み込むようにし、同時に(Applc)は(Windows)の(マルチメディア)【技術】の(サポート)を(強化)することになる。
74	1.1	2.2	今までは(Windows)用の(QuickTime)は(再生)しかできなかった。